

# Statistical Methods for Joint Analysis of Survival Time and Longitudinal Data

Jaeun Choi

A dissertation submitted to the faculty of the University of North Carolina at Chapel Hill in partial fulfillment of the requirements for the degree of Doctor of Philosophy in the Department of Biostatistics.

Chapel Hill  
2011

Approved by:

Advisor: Dr. Jianwen Cai

Advisor: Dr. Donglin Zeng

Reader: Dr. David Couper

Reader: Dr. Bahjat Qaqish

Reader: Dr. Andrew F. Olshan

© 2011  
Jaeun Choi  
ALL RIGHTS RESERVED

# Abstract

**JAEUN CHOI: Statistical Methods for Joint Analysis of Survival Time  
and Longitudinal Data.**

**(Under the direction of Dr. Jianwen Cai and Dr. Donglin Zeng.)**

In biomedical studies, researchers are often interested in the relationship between patients' characteristics or risk factors and both longitudinal outcomes such as quality of life measured over time and survival time. However, despite the progress in the joint analysis for longitudinal data and survival time, investigation on modeling approach to find which factor or treatment can simultaneously improve the patient's quality of life and reduce the risk of death has been limited. In this dissertation, we investigate joint modeling of longitudinal outcomes and survival time. We consider the generalized linear mixed models for the longitudinal outcomes to incorporate both continuous and categorical data and the stratified multiplicative proportional hazards model for the survival data. We study both Gaussian process and distribution free approaches for the random effect characterizing the joint process of longitudinal data and survival time.

We consider three estimation approaches in this dissertation. First, we consider the maximum likelihood approach with Gaussian process for random effects. The random effects, which are introduced into the simultaneous models to account for dependence between longitudinal outcomes and survival time due to unobserved factors, are assumed to follow a multivariate Gaussian process. The full likelihood, obtained by integrating the complete data likelihood over the random effects, is used for estimation. The Expectation-Maximization (EM) algorithm is used to compute the point

estimates for the model parameters, and the observed information matrix is adopted to estimate their asymptotic variances. Second, the normality assumption of random effects in the likelihood approach is relaxed. Assuming the underlying distribution of random effects to be unknown, we propose using a mixture of Gaussian distributions as an approximation in estimation. Weights of the mixture components are estimated with model parameters using the EM algorithm, and the observed information matrix is used for estimation of the asymptotic variances of the proposed estimators. For the two maximum likelihood approaches with and without normality assumption of random effects, asymptotic properties of the proposed estimators are investigated and their finite sample properties are assessed via simulation studies. Third, we consider a penalized likelihood approach. This approach is expected to be computationally less intensive than the maximum likelihood approach. It gives a penalty for regarding the random effect as a fixed effect in the likelihood and avoids the need to integrate the likelihood over random effects. The penalized likelihood is obtained through Laplace approximation. We compare the numerical performances of the penalized likelihood method and the EM algorithm used in maximum likelihood estimation for the simultaneous models with Gaussian process for random effects via simulation studies. All the proposed methods in this dissertation are illustrated with the real data from the Carolina Head and Neck Cancer Study (CHANCE).

# Acknowledgments

I would like to thank my dissertation advisors, Drs. Jianwen Cai and Donglin Zeng, for their expert guidance, deep insights, and thoughtful encouragement throughout my dissertation research process. The lessons I have learned under their direction and the experience I have gained during this process are invaluable. I am grateful to Dr. Jianwen Cai for her financial support as well.

I want to express my sincere gratitude to the committee members, Drs. David Couper, Bahjat Qaqish, and Andrew F. Olshan, for their constructive and perceptive comments. In particular, I am thankful to Dr. Andrew F. Olshan for providing the CHANCE data.

I would like to extend my heartfelt gratitude to the faculties and staffs of the Department of Biostatistics at UNC-CH for their help and encouragement in various ways during my study. I am also grateful to my former colleagues at the UNC Lineberger Comprehensive Cancer Center, where I had my first graduate research assistantship, for their warm consideration and motivation. I want to give my special thanks to all my friends and family for their love, support, and prayer, and to God who has always walked with me through the psalm 119:165, “Great peace have those who love your law, and nothing can make them stumble.”

# Table of Contents

<b>Abstract</b> . . . . .	iii
<b>List of Tables</b> . . . . .	x
<b>List of Figures</b> . . . . .	xii
<b>1 INTRODUCTION</b> . . . . .	1
1.1 Joint Analysis for Survival Time and Longitudinal Categorical Measure- ments of Quality of Life in Head and Neck Cancer Patients . . . . .	2
1.2 Joint Modeling of Survival Time and Longitudinal Outcomes with Flex- ible Random Effects . . . . .	3
1.3 Penalized Likelihood Approach for Joint Analysis of Survival Time and Longitudinal Outcomes . . . . .	4
<b>2 LITERATURE REVIEW</b> . . . . .	5
2.1 Failure Time Models . . . . .	5
2.1.1 Univariate failure time model . . . . .	6
2.1.2 Correlated failure time model . . . . .	8
2.2 Longitudinal Data Models and Methods . . . . .	10

2.2.1	Generalized linear model with random effects . . . . .	11
2.2.2	Maximum likelihood and conditional likelihood methods . . . . .	11
2.3	Joint Models of Failure Time and Longitudinal Data . . . . .	15
2.3.1	Failure time model with longitudinal covariates . . . . .	16
2.3.2	Simultaneous model of failure time and longitudinal data . . . . .	23
2.4	Penalized Quasi-Likelihood Approach . . . . .	27
2.4.1	Penalized quasi-likelihood in generalized linear mixed model . . . .	27
2.4.2	Bias correction in penalized quasi-likelihood . . . . .	31
<b>3</b>	<b>JOINT ANALYSIS FOR SURVIVAL TIME AND LONGITUDINAL CATEGORICAL MEASUREMENTS OF QUALITY OF LIFE IN HEAD AND NECK CANCER PATIENTS . . . . .</b>	<b>37</b>
3.1	Introduction . . . . .	37
3.2	The CHANCE Study . . . . .	41
3.3	Models and Inference Procedure . . . . .	44
3.3.1	Model formulation and notation . . . . .	44
3.3.2	Inference procedure . . . . .	46
3.3.3	EM algorithm – examples . . . . .	50
3.4	Asymptotic Properties . . . . .	52
3.5	Technical Details – Proofs for Asymptotic Properties . . . . .	54
3.5.1	Proof of consistency . . . . .	56
3.5.2	Proof of asymptotic normality . . . . .	72
3.5.3	Supplementary proofs . . . . .	81

3.6	Simulation Studies . . . . .	93
3.6.1	Binary longitudinal outcomes and survival time . . . . .	94
3.6.2	Poisson longitudinal outcomes and survival time . . . . .	95
3.7	Analysis of the CHANCE Study . . . . .	97
3.8	Concluding Remarks . . . . .	105
<b>4</b>	<b>JOINT MODELING OF SURVIVAL TIME AND LONGITUDINAL OUTCOMES WITH FLEXIBLE RANDOM EFFECTS . . . . .</b>	<b>107</b>
4.1	Introduction . . . . .	107
4.2	Models and Inference Procedure . . . . .	109
4.2.1	Model formulation and notation . . . . .	109
4.2.2	Inference procedure . . . . .	111
4.2.3	EM algorithm – examples . . . . .	115
4.3	Asymptotic Properties . . . . .	117
4.4	Technical Details – Proofs for Asymptotic Properties . . . . .	120
4.4.1	Proof of consistency . . . . .	122
4.4.2	Proof of asymptotic normality . . . . .	140
4.4.3	Supplementary proofs . . . . .	150
4.5	Simulation Studies . . . . .	165
4.5.1	Continuous longitudinal outcomes and survival time . . . . .	166
4.5.2	Binary longitudinal outcomes and survival time . . . . .	167
4.5.3	Sensitivity for model-misspecification . . . . .	169
4.5.4	Selection of the number of mixture distributions . . . . .	171



4.6	Analysis of the CHANCE Study . . . . .	175
4.7	Concluding Remarks . . . . .	180
<b>5</b>	<b>PENALIZED LIKELIHOOD APPROACH FOR JOINT ANALYSIS OF SURVIVAL TIME AND LONGITUDINAL OUTCOMES . . . .</b>	<b>184</b>
5.1	Introduction . . . . .	184
5.2	Model Formulation and Notation . . . . .	187
5.3	Estimation Procedure . . . . .	189
5.3.1	Laplace approximation . . . . .	191
5.3.2	Penalized likelihood . . . . .	194
5.3.3	Implementation . . . . .	196
5.4	Simulation Studies . . . . .	200
5.5	Analysis of the CHANCE Study . . . . .	206
5.6	Concluding Remarks . . . . .	211
<b>6</b>	<b>SUMMARY AND FUTURE RESEARCH . . . . .</b>	<b>213</b>
	<b>REFERENCES . . . . .</b>	<b>217</b>

# List of Tables

3.1	Descriptive statistics of predictors in the CHANCE study . . . . .	43
3.2	Descriptive statistics of outcome variables in the CHANCE study . . . .	44
3.3	Summary of simulation results of maximum likelihood estimation for binary longitudinal outcomes and survival time. . . . .	96
3.4	Summary of simulation results of maximum likelihood estimation for Poisson longitudinal outcomes and survival time. . . . .	98
3.5	Analyses results for the HNCS QoL and survival time of the CHANCE study . . . . .	102
4.1	Summary of simulation results of maximum likelihood estimation using mixtures of Gaussian distributions for random effects in the joint mod- eling of continuous longitudinal outcomes and survival time. . . . .	168
4.2	Summary of simulation results of maximum likelihood estimation using mixtures of Gaussian distributions for random effects in the joint mod- eling of binary longitudinal outcomes and survival time. . . . .	170
4.3	Summary of simulation results of sensitivity for model-misspecification .	172
4.4	Summary of simulation results: Frequencies on the selected number of Normal distributions in mixture (n=200) . . . . .	175
4.5	Results from final models of simultaneous and separate analyses for the Quality of Life and survival time for the CHANCE study . . . . .	178
5.1	Summary of simulation results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) in the si- multaneous modeling of binary longitudinal outcomes and survival time (n=200). . . . .	202

5.2	Summary of simulation results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) in the simultaneous modeling of binary longitudinal outcomes and survival time (n=400). . . . .	203
5.3	Analyses results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) for the Quality of Life and survival time for the CHANCE study . . . . .	210

# List of Figures

3.1	Estimated baseline cumulative hazards (solid line) with 95% confidence interval (dashed lines) by the simultaneous analysis of HNCS QoL longitudinal outcome and survival time . . . . .	103
3.2	Kaplan-Meier estimates (solid line) and the predicted survival probabilities based on the simultaneous analysis of HNCS QoL longitudinal outcome and survival time (dashed line) . . . . .	104
4.1	Density plots of random effects from simulation results of sensitivity for model-misspecification . . . . .	173
4.2	Relative bias plot of parameters in longitudinal and hazard models (thin and thick lines respectively) from simulation results of sensitivity for model-misspecification . . . . .	174
4.3	Estimated baseline cumulative hazards (solid line) with 95% confidence interval (dotted lines) by the simultaneous analysis of HNCS QoL longitudinal outcome and survival time . . . . .	181
4.4	The predicted conditional longitudinal trend based on the simultaneous models (solid line) and the empirical longitudinal trend (dotted line) based on the empirical longitudinal HNCS QoL satisfaction probabilities (dots) . . . . .	182
5.1	Plot of ratios of mean squared errors (MSEs) of maximum penalized likelihood estimator (MPLE) to maximum likelihood estimator (MLE) for parameters of predictors in longitudinal and hazard models (n=200, 400) . . . . .	207
5.2	Plot of ratios of user times of maximum penalized likelihood estimator (MPLE) to maximum likelihood estimator (MLE) (n=200, 400) . . . . .	208

# Chapter 1

## INTRODUCTION

The models for jointly analyzing the longitudinal data and survival time have been intensively developed in recent literature. Most models in such analysis would answer the questions regarding how one's quality of life affects time to death or given one's death time how quality of life changes over time. In many biomedical studies, it is of interest to assess the simultaneous effect of treatment or other factors on both patients' quality of life and risk of death taking into account the dependence between quality of life and survival time within a patient. To answer such questions, we consider the simultaneous modeling of longitudinal outcomes and survival time. In this dissertation, we propose two different maximum likelihood approaches with Gaussian process and without distributional assumption for the random effect. In addition, we consider penalized likelihood approach and compare its numerical performance with EM algorithm used in maximum likelihood estimation for the simultaneous models with Gaussian process for the random effect.

## 1.1 Joint Analysis for Survival Time and Longitudinal Categorical Measurements of Quality of Life in Head and Neck Cancer Patients

Patient survival and Quality of Life (QoL) are often recognized as two major outcome variables in the evaluation of head and neck cancer treatment in oncology community. QoL is important because it reflects the patients' critical physical, psycho-social, and emotional functions and it impacts communication with their caregivers. For the Carolina Head and Neck Cancer Study (CHANCE), we consider a joint analysis of survival time and longitudinal categorical QoL outcomes to find important variables for predicting both patients' QoL and risk of death. We first propose the maximum likelihood approach to simultaneously model the survival time with a stratified Cox proportional hazards model and longitudinal categorical outcomes with a generalized linear mixed model through random effects with normality assumption. Random effects, which are introduced into the simultaneous models to account for dependence between longitudinal outcomes and survival time due to unobserved factors, are assumed to follow a multivariate Gaussian process so that we can use the full likelihood for estimation by integrating the complete data likelihood over the random effects. EM algorithm is used to derive the point estimates for the model parameters, and the observed information matrix is adopted to estimate their asymptotic variances. The asymptotic properties of the proposed estimators are investigated and their finite sample properties are assessed via simulation studies. We illustrate the proposed approach with the real data of longitudinal Head and Neck Cancer Specific symptoms (HNCS) QoL and survival time from the CHANCE study.

## 1.2 Joint Modeling of Survival Time and Longitudinal Outcomes with Flexible Random Effects

In addition to the maximum likelihood approach with Gaussian process for random effects, we investigate a different method without any distributional assumption for random effects. Gaussian distribution is a convenient distribution often used for the random effects characterizing the joint process of longitudinal outcomes and survival time, and the likelihood approach relies heavily on the such normality assumption. However, this assumption may not be satisfied and the results could be misleading if the assumption is violated. These concerns motivate us to seek more robust estimation method which is not sensitive to the distributional assumption of random effects. Therefore, we relax the normality assumption of random effects by assuming the underlying distribution to be unknown. We propose to use a mixture of Gaussian distributions as an approximation in the estimation. Weights of the mixture components are estimated with model parameters using the EM algorithm, and the observed information matrix is used for the estimation of the asymptotic variances of the proposed estimators. The asymptotic properties of the proposed estimators are investigated and the method is demonstrated to perform well in finite samples via simulation studies. We also conduct simulation studies to examine the robustness of the mixture distribution. AIC and BIC criteria are adopted for selecting the number of mixtures, and the selection procedures are assessed through simulation studies. The proposed method is applied to the CHANCE study aforementioned.

### **1.3 Penalized Likelihood Approach for Joint Analysis of Survival Time and Longitudinal Outcomes**

We compare the numerical performances of the EM algorithm used in the maximum likelihood estimation with Gaussian process for random effect and another estimation method using the penalized likelihood. The penalized likelihood is expected to have less burden on computation because it treats the random effect as the fixed effect in the likelihood and penalized it. Consequently, no calculation is needed to integrate the likelihood over random effects. In SAS GLIMMIX procedure, penalized quasi-likelihood imposing the penalty in quasi-likelihood is already built and used for the analysis of the generalized linear mixed model. Accordingly, it is worthwhile to compare the numerical performances of the EM algorithm and the penalized likelihood method in maximum likelihood estimation. If the EM algorithm performs similarly to the penalized likelihood method on computational time, it will be better to maximize the full likelihood rather than the penalized likelihood. In the meantime, if the penalized likelihood method takes less time and provides unbiased and consistent estimates similar to those from EM algorithm, the penalized likelihood method will be preferred. We present the penalized likelihood obtained through Laplace approximation for our joint models and conduct simulation studies for performance comparison of penalized likelihood and EM algorithm. We also illustrate this comparison through the data analysis of the CHANCE study.



# Chapter 2

## LITERATURE REVIEW

In this section, we review the statistical literature for : 1) failure time models, 2) longitudinal data models and methods, 3) joint models of failure time and longitudinal data, and 4) penalized quasi-likelihood approach. The organization of the rest of this section is as following. We review literature on statistical methods for Cox proportional hazard models of univariate failure time and frailty models of correlated failure times in section 2.1, and, for generalized linear models with random effects and parameter estimation of longitudinal data in section 2.2. In section 2.3, we review the literature on statistical methods for joint models of failure time and longitudinal data. Lastly, we review penalized quasi-likelihood approach for generalized linear mixed model and bias correction for the estimator in section 2.4.

### 2.1 Failure Time Models

Failure time analysis or survival analysis addresses data of the form ‘time until an event occurs.’ The approaches were primarily developed in the medical and biological sciences, but are also broadly used in the social and economic sciences and engineering. A research question arising frequently in these areas is to determine whether or not certain variables are associated with the failure or survival times. There are two major reasons

why this problem cannot be handled via straightforward regression approaches: First, the dependent variable of interest (failure/survival time) is most likely not normally distributed, which is a serious violation of an assumption for ordinary least squares multiple regression. Survival times usually follow a skewed distribution. Second, there is the problem of censoring, that is, some observations will be incomplete.

We summarize the Cox proportional hazard model for the univariate failure time, which is not based on any assumptions concerning the nature or shape of the underlying survival distribution, in section 2.1.1, and the frailty model for the correlated failure time, which formulates the nature of dependence explicitly, as an extension of the Cox model in section 2.1.2.

### 2.1.1 Univariate failure time model

The Cox proportional hazards model (Cox, 1972) has been the most widely used procedure to study the effects of covariates on a failure time. The Cox model assumes that the hazard function for the failure time  $T$  associated with a covariate vector  $\mathbf{Z}$  is given by

$$\lambda(t|\mathbf{Z}) = \lambda_0(t) \exp\{\boldsymbol{\beta}_0^T \mathbf{Z}(t)\}, \quad t \geq 0 \quad (2.1)$$

where  $\lambda_0(t)$  is an unspecified baseline hazard function and  $\boldsymbol{\beta}_0$  is a  $p \times 1$  vector of unknown regression parameters. The model (2.1) is semi-parametric in that the effect of the covariates on the hazard is explicitly specified while the form of the baseline hazard function is unspecified. The model (2.1) assumes that hazard ratios are proportional across groups or subpopulations over time, and the regression coefficient  $\boldsymbol{\beta}_0$  represents the log hazard ratio for one unit increase in the corresponding covariate given that the other covariates in the model are held at the same value.

Let  $C$  denote the potential censoring time and  $X = \min(T, C)$  denote the observed time. Let  $N(t)$  denote the counting process,  $Y(t) = I(X \leq t)$  be an ‘at-risk’ indicator

process and  $\Delta = I(T \leq C)$  be an indicator for failure, where  $I(\cdot)$  is an indicator function. The failure time is assumed to be subject to independent right censorship. Let  $(T_i, C_i, \mathbf{Z}_i)(i = 1, \dots, n)$  be  $n$  independent replicates of  $(T, C, \mathbf{Z})$  and  $\tau$  denote the study end time.

The regression parameter  $\beta_0$  in (2.1) can be estimated by applying standard asymptotic likelihood procedure to the ‘partial’ likelihood function, introduced by Cox (1975),

$$L(\beta) = \prod_{i=1}^n \left[ \frac{\exp\{\beta^T \mathbf{Z}_i(T_i)\}}{\sum_{l=1}^n Y_l(T_i) \exp\{\beta^T \mathbf{Z}_l(T_i)\}} \right]^{\Delta_i},$$

where  $\mathbf{Z}_i(T_i)$  is the covariate vector for the subject failing at  $T_i$ , and  $\mathbf{Z}_l(T_i)$  is the corresponding covariate vector for the  $l$ -th member who is at risk at  $T_i$ . The estimator for  $\beta_0$ , denoted by  $\hat{\beta}$ , is obtained by the partial likelihood score function

$$\mathbf{U}(\beta) = \sum_{i=1}^n \Delta_i \left\{ \mathbf{Z}_i(X_i) - \frac{\mathbf{S}^{(1)}(\beta, X_i)}{\mathbf{S}^{(0)}(\beta, X_i)} \right\},$$

where  $\mathbf{S}^{(0)}(\beta, t) = n^{-1} \sum_{i=1}^n Y_i(t) \exp\{\beta' \mathbf{Z}_i(t)\}$ ,  $\mathbf{S}^{(1)}(\beta, t) = n^{-1} \sum_{i=1}^n Y_i(t) \exp\{\beta' \mathbf{Z}_i(t)\} \mathbf{Z}_i(t)$ . The maximum partial likelihood estimator  $\hat{\beta}$ , defined as the solution to the unbiased score equation  $\mathbf{U}(\beta) = \mathbf{0}$ , has been shown to be approximately normal in large samples with mean  $\beta_0$  and with a covariance matrix that can be consistently estimated by  $-\left\{ \frac{\partial \mathbf{U}(\beta)}{\partial \beta} \Big|_{\beta=\hat{\beta}} \right\}^{-1}$  (Andersen & Gill, 1982; Tsiatis, 1981). Iterative procedures, such as Newton-Raphson method and EM algorithm, are commonly used to solve the score equation.

Cox proportional hazards model has been extended from analyzing univariate failure time data to multivariate failure time data. Andersen & Gill (1982) and Fleming & Harrington (1991) extended Cox model in the expression of counting process which is more general and includes recurrent failures. Wei, Lin & Weissfeld (1989) and Hougaard (2000) extended Cox model to model multivariate failure times. As another extension

from Cox model, Clayton & Cuzick (1985) and Hougaard (2000) proposed frailty model for clustered failure data in which subjects may or may not experience the same type of event but they may be correlated because subjects are from the same cluster. In the frailty model, Cox proportional hazards model is used to model each individual's hazard function, and then an unobserved cluster-specific frailty is introduced into each model to account for within-cluster correlation. This frailty model is reviewed in the next section 2.1.2.

### **2.1.2 Correlated failure time model**

The Cox model (2.1) in the previous section 2.1.1 assumes the independent failure times. In many biomedical studies, however, the independence between failure times might be violated, which may arise because study subjects may be grouped in a manner that leads to dependencies within groups, or because individuals may experience multiple events. For such data, there are two main approaches: the marginal model approach which leaves the nature of dependence among related failure times completely unspecified and the frailty model approach which formulate the nature of dependence explicitly. When the interest resides in estimating the effect of risk factors and the correlation among the failure times are considered as a nuisance, the marginal model approach suits this purpose very well. However, in some settings, one might be interested in the strength and nature of dependencies among the failure time components, for which the frailty models have been proposed and studied by many authors. We focus on frailty model in this section.

The frailty model explicitly formulates the nature of the underlying dependence structure through an unobservable random variable. This unknown factor is usually called individual heterogeneity or frailty. The key assumption is that the failure times are conditionally independent given the value of the frailty. To illustrate this idea,

consider a Cox proportional hazards model for subject  $i$  with respect to the  $k$ th event :

$$\lambda_{ik}(t|w_i) = w_i \lambda_0(t) \exp\{\beta_0^T z_{ik}(t)\} \quad (2.2)$$

where the frailty terms  $\{w_i\}, i = 1, \dots, n$  are assumed to be independent and to arise from a common parametric density. The commonly used one is the gamma distribution, mostly for mathematical convenience. Various choices are possible for this density, which include the positive stable distributions, the inverse Gaussian distributions and the log-normal distributions. Note that  $\beta_0$  in (2.2) generally needs to be interpreted conditionally on the unobserved frailty. The frailty model approach is particularly sensible, when the strength of the dependence of failure times is of interest.

The parameter estimates are obtained through the EM algorithm, making use of the partial likelihood expression in the maximization step as shown in Klein (1992). An alternative approach is to use a penalized partial likelihood for the estimation of the shared frailty (Therneau & Grambsch, 2001).

Troxel & Esserman (2004) proposed a novel application of frailty models to assess the correlation between survival and quality of life in oncology. A frailty parameter is a random effect that allows the variability among clusters of measurements to be incorporated into survival models. The collected quality of life outcomes are dichotomized in order to apply the multivariate survival methods. In spite of the necessity of the conversion, the discretization of the quality of life scores from a continuous to a failure-time structure leads to the loss of information available from continuous quality of life data.

Ratcliffe, Guo, & TenHave (2004) proposed a joint model for the analysis of longitudinal and survival data in the presence of data clustering via a common frailty. While the existing models include subject-level random effects as the only random effects, two levels of nested random effects (subject-level random effects and cluster-level

frailty, with subjects nested within clusters) are used in the model with the responses linked at the higher cluster level. This additional level of random effects makes the model more flexible. They used a mixed effects model for the repeated measures that incorporates both subject- and cluster-level random effects, with subjects nested within clusters. A Cox frailty model is used for the survival model as it allows for between-cluster heterogeneity. Then they link the two responses via the common cluster-level random effects, or frailties, using a multivariate normal assumption for computation ease (Li & Lin, 2000). More joint models of survival and longitudinal data are reviewed using different models in section 2.3.

## 2.2 Longitudinal Data Models and Methods

The defining feature of a longitudinal study is that individuals are measured repeatedly through time. Longitudinal data require special statistical methods because the set of observations on one subject tends to be intercorrelated. This correlation must be taken into account to draw valid scientific inferences.

There are a variety of qualitatively different sources of random variation that might actually occur in practice and be included to construct the model for longitudinal data. The linear random effects model described in section 2.2.1 is one of three extensions of generalized linear models for longitudinal data: marginal, random effects, and transition models. The random effects model is most useful when the objective is to make inference about individuals rather than the population average. The parameter estimation in the generalized linear model with random effects can be carried out by both maximum likelihood and conditional likelihood approaches summarized in section 2.2.2.

### 2.2.1 Generalized linear model with random effects

The linear random effects model is applied where the response is assumed to be a linear function of explanatory variables with regression coefficients that vary from one individual to the next. This variability reflects natural heterogeneity due to unmeasured factors, which can be represented by a probability distribution. Correlation among observations for one person arises from their sharing unobservable variables,  $U_i$ .

The random effects GLM has the following general specifications:

1. Given the random effects  $\mathbf{U}_i$ , the responses  $Y_{i1}, \dots, Y_{in_i}$  are mutually independent and follows a distribution from the exponential family with density

$$f(y_{ij}|\mathbf{U}_i; \boldsymbol{\beta}) = \exp[\{(y_{ij}\theta_{ij} - \psi(\theta_{ij}))\}/\phi + c(y_{ij}, \phi)]. \quad (2.3)$$

The conditional moments,  $\mu_{ij} = E(Y_{ij}|\mathbf{U}_i) = \psi'(\theta_{ij})$  and  $v_{ij} = \text{Var}(Y_{ij}|\mathbf{U}_i) = \psi''(\theta_{ij})\phi$ , satisfy  $h(\mu_{ij}) = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{d}_{ij}^T \mathbf{U}_i$  and  $v_{ij} = v(\mu_{ij})\phi$  where  $h$  and  $v$  are known link and variance functions, respectively, and  $\mathbf{x}_{ij}$  and  $\mathbf{d}_{ij}$  are covariate vectors of length  $p$  and  $q$ , respectively.  $\mathbf{d}_{ij}$  is a subset of  $\mathbf{x}_{ij}$ .

2. The random effects,  $\mathbf{U}_i$ ,  $i = 1, \dots, m$ , are mutually independent and identically distributed with density function  $f(\mathbf{U}_i; G)$ .

Another fundamental assumption of the random effects model is that the  $\mathbf{U}_i$  are independent of the explanatory variables. A model of this type is sometimes referred to as a “latent variable” model (Bartholomew, 1987).

### 2.2.2 Maximum likelihood and conditional likelihood methods

Let  $\mathbf{U} = (U_1, \dots, U_m)$ . In maximum likelihood approach,  $\mathbf{U}$  is treated as a set of unobserved variables which is integrated out of the likelihood, adopting the assumption that the random effects follow a distribution such as Gaussian model with mean zero and

variance matrix  $\mathbf{G}$ . In conditional likelihood approach, the random effects is treated as if they were fixed parameters to be removed from the problem, so that we need not rely on the second assumption in the previous section 2.2.1.

Maximum likelihood approach treats  $\mathbf{U}_i$  as a sample of independent unobservable variables from a random effects distribution. Then, the likelihood function for the unknown parameter  $\boldsymbol{\delta}$ , which is defined to include both  $\boldsymbol{\beta}$  and the elements of  $\mathbf{G}$ , is

$$L(\boldsymbol{\delta}; \mathbf{y}) = \prod_{i=1}^m \int \prod_{j=1}^{n_i} f(y_{ij}|\mathbf{U}_i; \boldsymbol{\beta}) f(\mathbf{U}_i; \mathbf{G}) d\mathbf{U}_i, \quad (2.4)$$

which is the marginal distribution of  $\mathbf{Y}$  obtained by integrating the joint distribution of  $\mathbf{Y}$  and  $\mathbf{U}$  with respect to  $\mathbf{U}$ . In some special case such as the Gaussian linear model, the integral in (2.4) has a closed form, but for most non-Gaussian models, numerical methods are required for its evaluation.

To find the maximum likelihood estimate, we solve the score equations obtained by setting to zero the derivative with respect to  $\boldsymbol{\delta}$  of the log likelihood. Considering the ‘complete’ data for an individual to comprise  $(\mathbf{y}_i, \mathbf{U}_i)$  and restricting attention to canonical link functions (McCullagh & Nelder, 1989) for which  $\theta_{ij} = h(\mu_{ij}) = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{d}_{ij}^T \mathbf{U}_i$ , then the ‘complete data’ score function for  $\boldsymbol{\beta}$  has a particularly simple form

$$\mathbf{S}_{\boldsymbol{\beta}}(\boldsymbol{\delta}|\mathbf{y}, \mathbf{U}) = \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{x}_{ij} \{y_{ij} - \mu_{ij}(\mathbf{U}_i)\} = 0, \quad (2.5)$$

where  $\mu_{ij}(\mathbf{U}_i) = E(y_{ij}|\mathbf{U}_i) = h^{-1}(\mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{d}_{ij}^T \mathbf{U}_i)$ . The observed data score functions  $\mathbf{S}_{\boldsymbol{\beta}}(\boldsymbol{\delta}|\mathbf{y})$  are defined as the expectations of the complete data score functions  $\mathbf{S}_{\boldsymbol{\beta}}(\boldsymbol{\delta}|\mathbf{y}, \mathbf{U})$  in (2.5) with respect to the conditional distribution of  $\mathbf{U}$  given  $\mathbf{y}$ . This gives,

$$\mathbf{S}_{\boldsymbol{\beta}}(\boldsymbol{\delta}|\mathbf{y}) = \sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{x}_{ij} [y_{ij} - E\{\mu_{ij}(\mathbf{U}_i)|\mathbf{y}_i\}] = 0. \quad (2.6)$$



The score equations for  $\mathbf{G}$  can similarly be obtained as

$$\mathbf{S}_G(\boldsymbol{\delta}|\mathbf{y}) = \frac{1}{2}\mathbf{G}^{-1}\left\{\sum_{i=1}^m \mathbb{E}(\mathbf{U}_i\mathbf{U}_i^T|\mathbf{y}_i)\right\}\mathbf{G}^{-1} - \frac{m}{2}\mathbf{G}^{-1} = 0. \quad (2.7)$$

A common strategy to solve for the maximum likelihood estimate of  $\boldsymbol{\delta}$  is to use the EM algorithm (Dempster *et al.*, 1977). This algorithm iterates between an E-step, which involves evaluating the expectations in the above score equations (2.6) and (2.7) using the current values of the parameters, and an M-step, in which we solve the score equations to give updated parameter estimates. The dimension of the integration involved in the conditional expectation is  $q$ , the dimension of  $\mathbf{U}_i$ . When  $q$  is one or two, numerical integration techniques can be implemented reasonably easily. (e.g. Crouch & Spiegelman, 1990) For higher dimensional problems, Monte Carlo integration methods can be used. (e.g. the application of Gibbs sampling in Zeger & Karim, 1991)

Gaussian distribution is a convenient model used most for the random effects. When the regression coefficients are of primary interest, the specific form of the random effects distribution is less important. However, when the random effects are themselves the focus, inferences are more dependent on the assumptions about their distribution. Lange & Ryan (1989) suggested a graphical way to test the Gaussian assumption when the response variables are continuous. When the response variables are discrete, the same task becomes more difficult. Davidian & Gallant (1992) developed a non-parametric approach to estimate the random effects distribution with non-linear models.

In conditional likelihood approach for the generalized linear models with random effects (Diggle *et al.*, 1994; McCullagh & Nelder, 1989), the main idea is to treat the random effects  $\mathbf{U}_i$  as a set of nuisance parameters to be removed, and to estimate  $\boldsymbol{\beta}$  using the conditional likelihood of the data given the sufficient statistics for the  $\mathbf{U}_i$ .

Treating  $\mathbf{U}$  as fixed, the likelihood function for  $\boldsymbol{\beta}$  and  $\mathbf{U}$  is

$$L(\boldsymbol{\beta}, \mathbf{U}; \mathbf{y}) = \prod_{i=1}^m \prod_{j=1}^{n_i} f(y_{ij} | \boldsymbol{\beta}, \mathbf{U}_i) \propto \prod_{i=1}^m \prod_{j=1}^{n_i} \exp\{\theta_{ij} y_{ij} - \psi(\theta_{ij})\}, \quad (2.8)$$

where  $\theta_{ij} = \theta_{ij}(\boldsymbol{\beta}, \mathbf{U})$ . Restrict attention to canonical link functions (McCullagh & Nelder, 1989) for which  $\theta_{ij} = \mathbf{x}_{ij}^T \boldsymbol{\beta} + \mathbf{d}_{ij}^T \mathbf{U}_i$ , the likelihood in (2.8) can be written as

$$L(\boldsymbol{\beta}, \mathbf{U}; \mathbf{y}) = \exp\left\{\boldsymbol{\beta}^T \sum_{i,j} \mathbf{x}_{ij} y_{ij} + \sum_i \mathbf{U}_i^T \sum_j \mathbf{d}_{ij} y_{ij} - \sum_{i,j} \psi(\theta_{ij})\right\}.$$

Hence, the sufficient statistics for  $\boldsymbol{\beta}$  and  $\mathbf{U}_i$  are  $\sum_{i,j} \mathbf{x}_{ij} y_{ij}$  and  $\sum_{i,j} \mathbf{d}_{ij} y_{ij}$  respectively, and  $\sum_{i,j} \mathbf{d}_{ij} y_{ij}$  is sufficient for  $\mathbf{U}_i$  for fixed  $\boldsymbol{\beta}$ .

The conditional likelihood is proportional to the conditional distribution of the data given the sufficient statistics for the  $\mathbf{U}_i$ , and the contribution from subject  $i$  has the form

$$\begin{aligned} f\left(\mathbf{y}_i \mid \sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i; \boldsymbol{\beta}\right) &= \frac{f(\mathbf{y}_i; \boldsymbol{\beta}, \mathbf{U}_i)}{f(\sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i; \boldsymbol{\beta}, \mathbf{U}_i)} \\ &= \frac{f(\sum_j \mathbf{x}_{ij} y_{ij} = \mathbf{a}_i, \sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i; \boldsymbol{\beta}, \mathbf{U}_i)}{f(\sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i; \boldsymbol{\beta}, \mathbf{U}_i)}. \end{aligned} \quad (2.9)$$

For a discrete generalized linear model, this expression (2.9) can be written as

$$P\left(\mathbf{y}_i \mid \sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i; \boldsymbol{\beta}\right) = \frac{\sum_{R_{i1}} \exp(\boldsymbol{\beta}^T \mathbf{a}_i + \mathbf{U}_i^T \mathbf{b}_i)}{\sum_{R_{i2}} \exp(\boldsymbol{\beta}^T \sum_j \mathbf{x}_{ij} y_{ij} + \mathbf{U}_i^T \mathbf{b}_i)},$$

where  $R_{i1}$  is the set of possible values for  $\mathbf{y}_i$  such that  $\sum_j \mathbf{x}_{ij} y_{ij} = \mathbf{a}_i$  and  $\sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i$ , and  $R_{i2}$  is the set of values for  $\mathbf{y}_i$  such that  $\sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i$ . The conditional likelihood for  $\boldsymbol{\beta}$  given the data for all  $m$  individuals simplifies to

$$L(\boldsymbol{\beta} | \mathbf{y}, \sum_j \mathbf{d}_{ij} y_{ij} = \mathbf{b}_i) = \prod_{i=1}^m \frac{\sum_{R_{i1}} \exp(\boldsymbol{\beta}^T \mathbf{a}_i)}{\sum_{R_{i2}} \exp(\boldsymbol{\beta}^T \sum_{j=1}^{n_i} \mathbf{x}_{ij} y_{ij})}. \quad (2.10)$$

For simple cases such as the random intercept model, the conditional likelihood is reasonably easy to maximize (Breslow and Day, 1980). By the analogy with the usual score equations derived from the full likelihood, the score equations obtained from the conditional likelihood (2.10) can be used to get maximum conditional likelihood estimator for  $\beta$ .

The random effects generalized linear models in biostatistics have been studied enormously including the following literatures providing useful additional references: Laird & Ware (1982); Stiratelli *et al.*(1984); Gilmour *et al.*(1985); Schall (1990); Zeger & Karim (1991); Waclawiw & Liang (1993); Solomon & Cox (1992); Breslow & Clayton (1993); Drum & McCullagh (1993); Breslow & Lin (1995) and Lin & Breslow (1996).

## 2.3 Joint Models of Failure Time and Longitudinal Data

Joint analysis of survival time and repeated measurements has been intensively studied in recent literature. The most models which have been used in such analysis can be categorized into a selection model or a pattern mixture model. The selection model would answer the question regarding how one's quality of life affects death and the pattern-mixture model would describe the pattern of quality of life given one's death time. However, research interest is also often in finding which factor or treatment can simultaneously improve the patients' quality of life and reduce the risk of death, which can be studied by the simultaneous analysis of quality of life and survival.

Let  $\mathbf{Y}$  denote the longitudinal outcomes, for example, quality of life, then  $\mathbf{Y}$  are realizations of a latent process  $\tilde{\mathbf{Y}}$  measured with errors. Let  $T$  denote survival time.

A selection model focuses on estimating the distribution of  $T$  given  $\tilde{\mathbf{Y}}$ . Such a selection model has been studied by many authors: Tsiatis *et al.*(1995), Wulfsohn and Tsiatis (1997), Hu *et al.*(1998), Huang *et al.*(2001), and Xu and Zeger (2001a, b). Usually,  $\tilde{\mathbf{Y}}$  is modeled as a function of observed covariates and subject-specific random

effects; then it is fed into the model of  $T$  given  $\tilde{\mathbf{Y}}$  as a linear predictor. Selection model is reviewed in section 2.3.1.

In the pattern mixture model, a model is assumed for longitudinal outcome  $\mathbf{Y}$  conditional on survival time  $T$  (Wu and Carroll, 1988; Wu and Bailey, 1989; Hogan and Laird, 1997) and interest focuses on estimating parameters in the model for longitudinal outcome.

Simultaneous modeling serves the purpose to model both the process for quality of life,  $\mathbf{Y}$ , and survival time,  $T$ , given observed covariates  $\mathbf{X}$ . Zeng & Cai (2005) proposed such a model of quality of life  $\mathbf{Y}$  following normal distribution and survival time  $T$  by the observed covariates  $\mathbf{X}$  and by unobserved factors with normal density. This approach is reviewed in section 2.3.2. It is noted that this approach is different from either selection model or pattern-mixture model, although mathematically, all three models can be regarded as different ways of writing the distribution of  $(T, \mathbf{Y})$  given covariates.

### **2.3.1 Failure time model with longitudinal covariates**

Many longitudinal studies collect information on each participant both on a time-to-event and covariates which may vary with time. Recent interest has focused on joint models for longitudinal covariate data and a survival endpoint. A popular approach assumes that the longitudinal data follow a linear mixed effects model (Laird & Ware, 1982) and that survival depends on the covariate through a proportional hazards relationship with the underlying random effects. To implement the Cox model with time-dependent covariates, complete knowledge of the true covariate history for each subject is required; however, time-dependent covariates are generally measured intermittently, often at different times for each subject and with error. A naive approach is to substitute for each subject at each failure time in the Cox partial likelihood (Cox,

1975) the closest observed covariate value prior to that time, often termed ‘last value carried forward’. It is well known (Prentice, 1982) that substituting mis-measured values for true covariates in the Cox model leads to biased estimation. Another strategy for estimation of the proportional hazard regression parameters is a two-stage approach (Pawitan & Self, 1993; Tsiatis *et al.*, 1995) : First, the mixed effects model is fitted to data at each risk set assuming normality for both random effects and intra-subject error from which empirical Bayes estimates of the individual random effects are obtained as described by Laird and Ware (1982). Then, predictors for the covariate for each subject at each failure time based on the relevant fit are substituted for the true covariate values in the Cox partial likelihood. This approximate method uses regression calibration (Carroll *et al.*, 1995) to reduce bias of the naive approach but still yields biased estimators for large measurement error. Alternatively, the joint likelihood of the survival and longitudinal data may form the basis for inference. DeGruttola & Tu (1994) assumed the covariate process and survival times to be multivariate normal and fitted the model via parametric maximum likelihood. Wulfsohn & Tsiatis (1997) adopted the less rigid proportional hazards relationship and used nonparametric maximum likelihood, but continued to assumed normal random effects. Henderson *et al.*(2000) used normal random effects in Gaussian covariate stochastic processes. Faucett & Thomas (1996) assumed normality and took a Bayesian approach.

These strategies rely heavily on the assumption of normality of random effects characterizing the true covariate process; however, this assumption may be over-restrictive and the consequences if it is violated are unknown. Tsiatis & Davidian (2001) proposed a conditional score estimation for the proportional hazards model with longitudinal covariates with measurement-errors, which does not put any restrictions on the distribution of the random effects in covariate process by exploiting the conditional score approach of Stefanski & Carroll (1987).

In this section, we focus on the conditional score estimation approach by Tsiatis & Davidian (2001) since the fundamental idea of the maximum likelihood approach, which has been mostly studied with distributional assumption of random effects by many authors, is same as that reviewed in section 2.3.2.

For each subject  $i$  ( $i = 1, \dots, n$ ), let  $T_i$  and  $C_i$  denote time to failure and censoring, respectively, where time on study  $V_i = \min(T_i, C_i)$  and failure indicator  $\Delta_i = I(T_i \leq C_i)$  are observed; all variables are independent across  $i$ . Let  $Z_i$  denote time-independent covariates and  $X_i(u)$  denote time-dependent covariates at time  $u$  for subject  $i$ ; for simplicity, assume  $X_i(u)$  is scalar, but generalization to vector-valued  $X_i(u)$  is straightforward. Assume that  $X_i(u)$  follows a subject-specific linear model  $X_i(u) = \alpha_{0i} + \alpha_{1i}u$ , where  $\alpha_i = (\alpha_{0i}, \alpha_{1i})^T$  are the intercept and slope for  $i$ . The covariate process  $X_i(u)$  is not directly observed; rather, longitudinal measurements  $W_i(t_{ij})$  are obtained at ordered times  $t_i = (t_{i1}, \dots, t_{im_i})^T$ , for  $t_{im_i} \leq V_i$ , where  $W_i(t_{ij}) = X_i(t_{ij}) + e_{ij}$ , with  $e_i = (e_{i1}, \dots, e_{im_i})^T$ . The errors  $e_{ij}$  reflect uncertainty in measuring  $X_i(u)$  at  $t_{ij}$  and are assumed identically normally distributed and independent with mean zero and variance  $\sigma^2$ , independent of  $(T_i, C_i, \alpha_i, Z_i, t_i, m_i)$ . More precisely,

$$(e_i | T_i, C_i, \alpha_i, Z_i, t_i, m_i) \sim N_{m_i}(0, \sigma^2 I_{m_i}),$$

where  $I_{m_i}$  the  $m_i$ -dimensional identity matrix.

The survival model assumes that the hazard of failure is related to  $X_i(u)$  and  $Z_i$  through a proportional hazards regression model; that is,

$$\begin{aligned} \lambda_i(u) &= \lim_{du \rightarrow 0} du^{-1} \text{pr}\{u \leq T_i < u + du | T_i \geq u, \alpha_i, Z_i, C_i, e_i(u), t_i(u)\} \\ &= \lim_{du \rightarrow 0} du^{-1} \text{pr}\{u \leq T_i < u + du | T_i \geq u, \alpha_i, Z_i\} \\ &= \lambda_0(u) \exp\{\gamma X_i(u) + \eta^T Z_i\}, \end{aligned} \tag{2.11}$$

where  $\lambda_0(u)$  denotes an unspecified baseline hazard function, the collection of times of longitudinal measurements up to and including  $u$  is denoted by  $t_i(u) = (t_{ij} \leq u)$ ,  $e_i(u) = (e_{ij} : t_{ij} \leq u)$ , and  $\eta$  is  $(q \times 1)$ . The model (2.11) shows explicitly the nature of the assumption that timing of measurements and censoring are noninformative. Interest focuses on estimation of the parameters  $\gamma$  and  $\eta$ .

Let  $\hat{X}_i(u)$  be the ordinary least squares estimator of  $X_i(u)$  using all the longitudinal data up to and including time  $u$ , that is based on  $t_i(u)$ . This requires at least two longitudinal measurements on  $i$  up to and including  $u$ , for  $t_{i2} \leq u$ . Define the counting process increment

$$dN_i(u) = I(u \leq V_i < u + du, \Delta_i = 1, t_{i2} \leq u)$$

and the 'at risk' process  $Y_i(u) = I(V_i \geq u, t_{i2} \leq u)$ ; that is,  $dN_i(u)$  puts point mass at time  $u$  corresponding to the observed death time for the  $i$ -th subject as long as this occurs after the second longitudinal measurement, and  $Y_i(u)$  is the indicator that subject  $i$  is at risk with at least two longitudinal measurement at time  $u$ . Then the estimator  $\hat{X}_i(u)$ , conditional on  $\{\alpha_i, t_i(u), Y_i(u) = 1, Z_i\}$ , is normally distributed with mean  $X_i(u) = \alpha_{0i} + \alpha_{1i}u$  and variance  $\sigma^2\theta_i(u)$ , the usual variance of the estimated mean  $\hat{X}_i(u)$  at  $u$  using data up to and including  $u$ , which depends on timing of measurements for  $i$  up to and including  $u$ . For  $X_i(u) = \alpha_{0i} + \alpha_{1i}u$ ,  $\theta_i(u) = 1/m_{i,u} + (u - \bar{t}_{i,u})^2/SS_{i,u}$ , where  $t_i(u)$  contains  $m_{i,u}$  time-points  $t_{ij}$  with mean  $\bar{t}_{i,u}$ ,  $SS_{i,u} = \sum_{j=1}^{m_{i,u}} (t_{ij} - \bar{t}_{i,u})^2$ .

At any time  $u$ , given that  $i$  is at risk at time  $u$  so that  $Y_i(u) = 1$ , random effects  $\alpha_i$ , longitudinal measurements taken up to and including time  $u$  at times  $t_i(u)$ , and time-independent covariates  $Z_i$ , the conditional density for  $\{dN_i(u) = r, \hat{X}_i(u) = x\}$  is

$$\text{pr}\{dN_i(u) = r | Y_i(u) = 1, \hat{X}_i(u) = x, \alpha_i, Z_i, t_i(u)\} \times \text{pr}\{\hat{X}_i(u) = x | Y_i(u) = 1, \alpha_i, Z_i, t_i(u)\},$$

which equals

$$\frac{[\lambda_0(u)du \exp\{\gamma X_i(u) + \eta^T Z_i\}]^r [1 - \lambda_0(u)du \exp\{\gamma X_i(u) + \eta^T Z_i\}]^{1-r}}{\{2\pi\sigma^2\theta_i(u)\}^{\frac{1}{2}}} \exp\left[-\frac{\{x - X_i(u)\}^2}{2\sigma^2\theta_i(u)}\right];$$

thus, the conditional likelihood of  $\{dN_i(u), \hat{X}_i(u)\}$  given  $\{Y_i(u) = 1, \alpha_i, Z_i, t_i(u)\}$ , up to order  $du$ , is

$$\begin{aligned} & [\lambda_0(u)du \exp\{\gamma X_i(u) + \eta^T Z_i\}]^{dN_i(u)} \frac{\exp[-\{\hat{X}_i(u) - X_i(u)\}^2 / \{2\sigma^2\theta_i(u)\}]}{\{2\pi\sigma^2\theta_i(u)\}^{\frac{1}{2}}} \\ &= \exp\left[X_i(u) \left\{ \gamma dN_i(u) + \frac{\hat{X}_i(u)}{\sigma^2\theta_i(u)} \right\}\right] \frac{\{\lambda_0(u) \exp(\eta^T Z_i) du\}^{dN_i(u)}}{\{2\pi\sigma^2\theta_i(u)\}^{\frac{1}{2}}} \exp\left\{-\frac{\hat{X}_i^2(u) + X_i^2(u)}{2\sigma^2\theta_i(u)}\right\}. \end{aligned}$$

This representation implies that, conditional on  $Y_i(u) = 1$ ,

$$S_i(u, \gamma, \sigma^2) = \gamma\sigma^2\theta_i(u)dN_i(u) + \hat{X}_i(u)$$

is a complete sufficient statistic for  $\alpha_i$ , suggesting that, at each time  $u$ , conditioning on  $S_i(u, \gamma, \sigma^2)$  would remove the dependence of the conditional distribution on the random effects  $\alpha_i$ . Then, the conditional intensity process defined as

$$\lim_{du \rightarrow 0} du^{-1} \text{pr}\{dN_i(u) = 1 | S_i(u, \gamma, \sigma^2), Z_i, t_i(u), Y_i(u)\}$$

is equal to  $\lambda_0(u) \exp\{\gamma S_i(u, \gamma, \sigma^2) - \gamma^2\sigma^2\theta_i(u)/2 + \eta^T Z_i\} Y_i(u)$ . Reasoning underlying the conditional score estimator follows by analogy with that for estimators for the proportional hazards model with no measurement error.

The conditional intensity of  $dN(u) = \sum_{j=1}^n dN_j(u)$ , given  $\{S_i(u, \gamma, \sigma^2), Z_i, t_i(u), Y_i(u), i = 1, \dots, n\}$ , is  $\lambda_0(u) E_0(u, \gamma, \eta, \sigma^2)$ , where  $E_0(u, \gamma, \eta, \sigma^2) = \sum_{j=1}^n E_{0j}(u, \gamma, \eta, \sigma^2)$ ,

$$E_{0j}(u, \gamma, \eta, \sigma^2) = \exp\{\gamma S_j(u, \gamma, \sigma^2) - \gamma^2\sigma^2\theta_j(u)/2 + \eta^T Z_j\} Y_j(u).$$



This suggests that a reasonable estimator for  $d\Lambda_0(u) = \lambda_0(u)du$  is given by

$$d\hat{\Lambda}_0(u) = dN(u)/E_0(u, \gamma, \eta, \sigma^2).$$

By analogy with the usual score equations derived from the partial likelihood in a proportional hazard model,  $(\gamma, \eta)$  can be obtained by solving the  $(q+1) \times 1$  set of estimating equations

$$\sum_{i=1}^n \int \{S_i(u, \gamma, \sigma^2), Z_i^T\}^T \{dN_i(u) - E_{0i}(u, \gamma, \eta, \sigma^2)d\hat{\Lambda}_0(u)\} = 0,$$

which upon substitution of  $d\hat{\Lambda}_0(u)$  for  $d\Lambda_0(u)$ , may be written as

$$\sum_{i=1}^n \int \{S_i(u, \gamma, \sigma^2), Z_i^T\}^T \left\{ dN_i(u) - \frac{dN(u)E_{0i}(u, \gamma, \eta, \sigma^2)}{E_0(u, \gamma, \eta, \sigma^2)} \right\} = 0 \quad (2.12)$$

Defining  $E_{1j}(u, \gamma, \eta, \sigma^2) = \{S_j(u, \gamma, \sigma^2), Z_j^T\}^T \exp\{\gamma S_j(u, \gamma, \sigma^2) - \gamma^2 \sigma^2 \theta_j(u)/2 + \eta^T Z_j\} Y_j(u)$ ,

$$E_1(u, \gamma, \eta, \sigma^2) = \sum_{j=1}^n E_{1j}(u, \gamma, \eta, \sigma^2),$$

and interchanging the sums in (2.12), the estimating equations are expressed as

$$\sum_{i=1}^n \int \left[ \{S_i(u, \gamma, \sigma^2), Z_i^T\}^T - \frac{E_1(u, \gamma, \eta, \sigma^2)}{E_0(u, \gamma, \eta, \sigma^2)} \right] dN_i(u) = 0. \quad (2.13)$$

With no measurement error,  $\sigma^2 = 0$ , (2.13) is identical to the score equations for the maximum partial likelihood estimator of Cox (1975). With  $X_i(u)$  time-independent and  $\sigma^2$  known, the equations are asymptotically equivalent to those proposed by Nakamura (1992). There is an alternative semiparametric estimator with time-independent covariates studied by Buzas (1998).

There are more recent literatures on the selection model. Ribaud, Thompson &

Allen-Mersh (2000) proposed the application of a random effect selection model in the form of a trivariate Normal model for the joint analysis of QoL response and log survival time. The trivariate Normal model presented by Schluchter is a model that has been discussed in the context of drop-out. This model is a random effect selection model that assumes that the random parameters of a subject's underlying response profile such as intercept and slope of QoL response over time, and the logarithm of the survival time follow a trivariate Normal distribution.

Xu & Zeger (2001) developed latent variable models for joint analysis of longitudinal data comprising repeated measures and times to events, starting with the latent variable formulation of Fawcett and Thomas(1996), and extending and adapting it to the problem of identifying whether a longitudinal variable  $Y$  is a useful auxiliary or surrogate variable for event time  $T$  given other covariates. The linking linear predictor of  $Y$  and  $T$  was assumed to follow a Gaussian stochastic process suggested by Diggle (1988).

Song, Davidian, & Tsiatis (2002) assumed that the random effects have distribution in a plausible class with smooth densities, in mixed effects model for longitudinal covariates process belonging to proportional hazards model of event time. They used a class of smooth densities studied by Gallant & Nychka (1987). One speculation of Song *et al.*(2002) is that it is possible that the likelihood based approach using normality yields consistent estimator even when normality is a mis-specification under certain 'nice' conditions through their simulations.

Zeng & Cai (2005) provided the rigorous proof of the consistency of the maximum likelihood estimators and derivation of their asymptotic distributions because there was a lack of theoretical justification of the asymptotic properties for the MLEs even if the ML estimation has been extensively used in the joint analysis of repeated measurements and survival time and has been shown to perform well in numerical studies (Hu, Tsi-

atis & Davidian 1998). Their theoretical results further confirmed that nonparametric maximum likelihood estimation, which was proposed in the literature (Wu & Carroll, 1988; Tsiatis, DeGruttola and Wulfsohn, 1995; Wulfsohn & Tsiatis, 1997), provided efficient estimation. Additionally, it was also shown that the profile likelihood function can be used to give a consistent estimator for the asymptotic variance of the regression coefficients.

Tseng, Hsieh, & Wang (2005) proposed the joint modeling of longitudinal covariates and survival time using accelerated failure time since the accelerated failure time model is an attractive alternative to the Cox model when the proportionality assumption is not appropriate to describe the relationship between the survival time and longitudinal covariates. Hsieh, Tseng, & Wang (2006) recently studied maximum likelihood approach for the joint modelling of survival time and longitudinal covariates in details more.

Song & Wang (2007) proposed semiparametric approaches for joint modeling of longitudinal covariates and survival data with time-varying coefficients. To deal with covariate measurement error, they proposed a local corrected score estimator and a local conditional score estimator which are semiparametric methods in the sense that there is no distributional assumption needed for the underlying true covariates. Li, Wang, & Wang (2007) proposed score functions, named generalized sufficient and conditional scores, for the joint models of a primary endpoint and multiple longitudinal covariate processes by adjusting the bias resulted from the approaches by Li, Zhang & Davidian (2004).

### **2.3.2 Simultaneous model of failure time and longitudinal data**

In many biomedical studies, it is often of interest to investigate the simultaneous effect of treatment or other factors on both patients' quality of life and risk of death. Xu & Zeger

(2001b) and Zeng & Cai (2005) proposed similar simultaneous models of continuous longitudinal outcome  $\mathbf{Y}$  and survival time  $T$ . However, while in the model by Xu & Zeger a common latent process is shared by both  $\mathbf{Y}$  and  $T$ , Zeng & Cai allow individual random effects to affect quality of life and survival time very differently.

In the approach by Zeng & Cai (2005), quality of life and survival time are modeled through parametric and semiparametric models, respectively, assuming a linear mixed effect model for the longitudinal outcomes of quality of life and a multiplicative hazard model for survival time. In both models, observed covariates, which are included as predictors, are assumed to be either time-independent or external time-dependent variables. Unobserved factors enter the models as subject-specific random effects so as to account for unobserved heterogeneity.

For subject  $i$  given  $T > t$  and the observed history till time  $t$ , the longitudinal outcome of quality of life  $Y_i(t)$  at time  $t$  follows the linear mixed effect model,

$$Y_i(t) = \mathbf{X}_i(t)\boldsymbol{\beta} + \tilde{\mathbf{X}}_i(t)\mathbf{a}_i + \epsilon_i(t),$$

where  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  are the row vectors of the observed covariates and can be completely different or share some components,  $\epsilon_i(t)$  is a white noise process with mean zero and variance  $\sigma_y^2$ , and  $\mathbf{a}_i$  denotes a vector of subject-specific random effect of dimension  $k_0$  following a multivariate normal distribution with mean zero and covariance matrix  $\boldsymbol{\Sigma}_a$ , and  $\boldsymbol{\beta}$  is a column vector of coefficients for  $\mathbf{X}_i(t)$ . The random effect  $\mathbf{a}_i$  reflects the unobserved heterogeneity and is allowed to differ for different levels of covariates  $\tilde{\mathbf{X}}_i(t)$ .

For the survival time  $T_i$  given the observed covariates, the observed history till time  $t$ , and random effect  $\mathbf{a}_i$ , the conditional hazard rate function is assumed to follow a

multiplicative hazards model,

$$\lambda(t) \exp\{\tilde{\mathbf{W}}_i(t)(\boldsymbol{\phi} \circ \mathbf{a}_i) + \mathbf{W}_i(t)\boldsymbol{\gamma}\},$$

where  $\mathbf{W}_i(t)$  and  $\tilde{\mathbf{W}}_i(t)$  are the row vectors of the observed covariates and may share the same components,  $\boldsymbol{\phi}$  is a vector of parameters, and  $\lambda(t)$  is the baseline hazard rate function, and  $\boldsymbol{\gamma}$  is a column vector of coefficients for  $\mathbf{W}_i(t)$ . For the dependence parameter  $\boldsymbol{\phi}$  between quality of life and survival time,  $\boldsymbol{\phi} = 0$  means the dependence can be fully attributed to the observed covariates, and  $\boldsymbol{\phi} \neq 0$  implies that such dependence may also be due to some latent variables.

Supposing the survival time is possibly right censored with completely random right-censored time  $C_i$ , and assuming  $N_i$ , the number of the observed quality of life measurements for subject  $i$ , to be non-informative about parameters of interest, the observed data from  $n$  subjects are

$$(N_i, Y_i^j, \mathbf{X}_i^j, \tilde{\mathbf{X}}_i^j), \quad j = 1, \dots, N_i, i = 1, \dots, n,$$

$$(Z_i, \Delta_i, \{(\mathbf{W}_i(t), \tilde{\mathbf{W}}_i(t)) : t \leq Z_i\}), \quad i = 1, \dots, n,$$

where for subject  $i$ ,  $(Y_i^j, \mathbf{X}_i^j, \tilde{\mathbf{X}}_i^j)$  is the  $j$ -th observation of  $(\mathbf{Y}_i, \mathbf{X}_i, \tilde{\mathbf{X}}_i)$ ,  $Z_i = \min(T_i, C_i)$ , and  $\Delta_i = I(T_i \leq C_i)$ . Interests are estimating and making inference on the parameters  $\boldsymbol{\theta} = (\sigma_y, \boldsymbol{\Sigma}_a, \boldsymbol{\beta}, \boldsymbol{\phi}, \boldsymbol{\gamma})$  and the baseline cumulative hazard function  $\Lambda(t) = \int_0^t \lambda(s)ds$ .

Their estimation approach is likelihood-based. In the maximum likelihood method, given the random effects for the  $i$ -th subject, the observed quality of life with a multivariate Gaussian distribution is independent of the observed survival time with proportional hazard assumption, and the likelihood contribution of the  $i$ -th subject is integrated over the random effects in the joint models. Then, the observed likelihood

function for  $(\boldsymbol{\theta}, \Lambda)$  is expressed as

$$\begin{aligned}
L = & \prod_{i=1}^n \int_{\mathbf{a}} \left[ (2\pi\sigma_y^2)^{-N_i/2} \exp\{-(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta} - \tilde{\mathbf{X}}_i\mathbf{a})^T(\mathbf{Y}_i - \mathbf{X}_i\boldsymbol{\beta} - \tilde{\mathbf{X}}_i\mathbf{a})/2\sigma_y^2\} \right. \\
& \lambda(Z_i)^{\Delta_i} \exp\left\{ \Delta_i(\tilde{\mathbf{W}}_i(Z_i)(\boldsymbol{\phi} \circ \mathbf{a}) + \mathbf{W}_i(Z_i)\gamma) - \int_0^{Z_i} e^{\tilde{\mathbf{W}}_i(s)(\boldsymbol{\phi} \circ \mathbf{a}) + \mathbf{W}_i(s)\gamma} d\Lambda(s) \right\} \\
& \left. (2\pi)^{-k_0/2} |\boldsymbol{\Sigma}_a|^{-1/2} \exp\{-\mathbf{a}^T \boldsymbol{\Sigma}_a^{-1} \mathbf{a}/2\} \right] d\mathbf{a},
\end{aligned}$$

where  $\mathbf{Y}_i$  denotes the vector of  $(Y_i^1, \dots, Y_i^{N_i})^T$ ,  $\mathbf{X}_i$  denotes the matrix of  $((\mathbf{X}_i^1)^T, \dots, (\mathbf{X}_i^{N_i})^T)^T$ ,  $\tilde{\mathbf{X}}_i$  denotes  $((\tilde{\mathbf{X}}_i^1)^T, \dots, (\tilde{\mathbf{X}}_i^{N_i})^T)^T$ , and  $k_0$  is the dimension of  $\mathbf{a}$ .

EM algorithms are employed for the maximum likelihood estimates for  $(\boldsymbol{\theta}, \Lambda)$  over a set in which  $\boldsymbol{\theta}$  is in a bounded set and  $\Lambda$  belongs to a space consisting of all the increasing functions with  $\Lambda(0) = 0$ . It is clear that the maximum likelihood estimate for  $\Lambda$  can be chosen as a step function with jumps only at the observed failure times. In the EM algorithm,  $\mathbf{a}_i$  is considered as the missing statistics for  $i = 1, \dots, n$ . Therefore, the M-step solves the conditional score equation from the complete data given the observations, where the conditional expectation can be evaluated in the E-step. The iteration between E-step and M-step is conducted until the estimates converge. The final maximum likelihood estimate for  $(\boldsymbol{\theta}, \Lambda)$  is denoted by  $(\hat{\boldsymbol{\theta}}, \hat{\Lambda})$ .

The variance estimator for  $\hat{\boldsymbol{\theta}}$  is obtained by using the profile likelihood function whose logarithm is defined as  $pl_n(\boldsymbol{\theta}) = \max_{\Lambda} n^{-1} \sum_{i=1}^n q_i(\boldsymbol{\theta}, \Lambda)$  where  $q_i(\boldsymbol{\theta}, \Lambda)$ ,  $i = 1, \dots, n$ , is the logarithm of the observed likelihood function for the  $i$ -th subject. Particularly, an efficient algorithm, which is based on the EM-algorithm to calculate the profile likelihood function, is proposed and called as the PEME algorithm (partial expectation, maximization and evaluation).

## 2.4 Penalized Quasi-Likelihood Approach

In the view of the cumbersome and often intractable numerical integrations required for a full likelihood analysis, several suggestions were made for approximate inference in generalized linear mixed models and other nonlinear variance component models. One approach was proposed by Breslow & Clayton (1993) with some modifications to a Laplace expansion in order to motivate standard estimating equations that may be solved by iterative application of normal theory variance components procedures. In this section, we mainly review the penalized quasi-likelihood for generalized linear mixed model proposed by Breslow & Clayton (1993), and the bias correction in the penalized quasi-likelihood estimators proposed by Breslow & Lin (1995).

### 2.4.1 Penalized quasi-likelihood in generalized linear mixed model

The penalized quasi-likelihood (PQL) method exploited by Green (1987) for semiparametric regression analysis is available for inference in hierarchical models where the focus is on shrinkage estimation of the random effects (Robinson 1991). The PQL was proposed as an approximate Bayes procedure for some commonly occurring GLMM's by Laird (1978). Breslow & Clayton (1993) considered two closely related approximate methods (Penalized Quasi-Likelihood and Marginal Quasi-Likelihood) of inference in GLMM's and investigated their suitability for practical work by means of Monte Carlo studies and illustrative applications. Here we review only the PQL in their paper. They provided the PQL criterion motivated by approximating the integrated quasi-likelihood, and developed an approximate GLM for the marginal distribution of the data. The approximate GLM is related to the generalized estimating equation approach of Zeger et al. (1988). They note that PQL tends to underestimate somewhat the variance com-

ponents and (in absolute value) fixed effects when applied to clustered binary data, but the situation improves rapidly for binomial observations having denominators greater than one.

Within the framework of the generalized linear mixed model (GLMM), given an unobserved vector of random effects, observations are assumed to be conditionally independent with means that depend on the linear predictor through a specified link function and conditional variances that are specified by a variance function, known prior weights and a scale factor. The random effects are assumed to be normally distributed with mean zero and dispersion matrix depending on unknown variance components.

Consider hierarchical model and denote  $y_i$ ,  $i = 1, \dots, n$ , as the  $i$ -th observation of a univariate response variable with two vectors  $\mathbf{x}_i$  and  $\mathbf{z}_i$  of explanatory variables associated with the fixed and random effects respectively. The  $n$  responses may be blocked in some way, for example when they involve repeated measures on the same subject. Suppose that, given a  $q$ -dimensional vector  $\mathbf{b}$  of random effects, the  $y_i$  are conditionally independent with means  $E(y_i|\mathbf{b}) = \mu_i^b$  and variances  $\text{Var}(y_i|\mathbf{b}) = \phi a_i v(\mu_i^b)$ , where  $v(\cdot)$  is a specified variance function,  $a_i$  is a known constant (e.g., the reciprocal of a binomial denominator) and  $\phi$  is a dispersion parameter that may or may not be known. The conditional mean is related to the linear predictor  $\eta_i^b = \mathbf{x}_i^T \boldsymbol{\alpha} + \mathbf{z}_i^T \mathbf{b}$  by the link function  $g(\mu_i^b) = \eta_i^b$ , with inverse  $h = g^{-1}$ , where  $\boldsymbol{\alpha}$  is a  $p$  vector of fixed effects. Denoting the observation vector by  $\mathbf{y} = (y_1, \dots, y_n)^T$  and the design matrices with rows  $\mathbf{x}_i^T$  and  $\mathbf{z}_i^T$  by  $\mathbf{X}$  and  $\mathbf{Z}$ , the conditional mean satisfies

$$E(\mathbf{y}|\mathbf{b}) = h(\mathbf{X}\boldsymbol{\alpha} + \mathbf{Z}\mathbf{b}).$$

Assume that  $\mathbf{b}$  has a multivariate normal distribution with mean  $\mathbf{0}$  and covariance matrix  $\mathbf{D} = \mathbf{D}(\boldsymbol{\theta})$  depending on an unknown vector  $\boldsymbol{\theta}$  of variance components. In binomial, Poisson, and hypergeometric specifications for the conditional distribution



of  $y_i$ , the dispersion parameter  $\phi$  is fixed at unity. In other cases, however, it may be estimated together with  $\boldsymbol{\theta}$  as a parameter in the covariance matrix of the marginal distribution of  $\mathbf{y}$ .

The integrated quasi-likelihood function used to estimate  $(\boldsymbol{\alpha}, \boldsymbol{\theta})$  is defined by

$$e^{ql(\boldsymbol{\alpha}, \boldsymbol{\theta})} \propto |\mathbf{D}|^{-1/2} \int \exp \left[ -\frac{1}{2\phi} \sum_{i=1}^n d_i(y_i; \mu_i^b) - \frac{1}{2} \mathbf{b}^T \mathbf{D}^{-1} \mathbf{b} \right] d\mathbf{b}, \quad (2.14)$$

where  $d_i(y, \mu) = -2 \int_y^\mu \frac{y-u}{a_i v(u)} du$  denotes the deviance measure of fit. If, conditionally on  $\mathbf{b}$ , the observations are drawn from a linear exponential family with variance function  $v(\cdot)$ , then the deviance is well known to equal to the scaled difference  $2\phi\{l(y; y, \phi) - l(y; \mu, \phi)\}$ , where  $l(y; \mu, \phi)$  denotes the conditional likelihood of  $y$  given its mean  $\mu$  (McCullagh & Nelder 1989). In this case  $ql(\boldsymbol{\alpha}, \boldsymbol{\theta})$  represents the true log-likelihood of the data. The primary difficulty in implementing full likelihood inference lies in the integrations needed to evaluate  $ql$  and its partial derivatives.

The equation (2.14) can be written as  $c|\mathbf{D}|^{-1/2} \int e^{-\boldsymbol{\kappa}(\mathbf{b})} d\mathbf{b}$ , and then applied with Laplace's method for integral approximation (Barndorff-Nielsen & Cox 1989; Tierney & Kadane 1986). Let  $\boldsymbol{\kappa}'$  and  $\boldsymbol{\kappa}''$  denote the  $q$  vector and  $q \times q$  dimensional matrix of first- and second-order partial derivatives of  $\boldsymbol{\kappa}$  with respect to  $\mathbf{b}$ . Ignoring the multiplicative constant  $c$ , the approximation yields

$$ql(\boldsymbol{\alpha}, \boldsymbol{\theta}) \approx -\frac{1}{2} \log |\mathbf{D}| - \frac{1}{2} \log |\boldsymbol{\kappa}''(\tilde{\mathbf{b}})| - \boldsymbol{\kappa}(\tilde{\mathbf{b}}), \quad (2.15)$$

where  $\tilde{\mathbf{b}} = \tilde{\mathbf{b}}(\boldsymbol{\alpha}, \boldsymbol{\theta})$  denotes the solution to

$$\boldsymbol{\kappa}' = -\sum_{i=1}^n \frac{(y_i - \mu_i^b) \mathbf{z}_i}{\phi a_i v(\mu_i^b) g'(\mu_i^b)} + \mathbf{D}^{-1} \mathbf{b} = 0$$

that minimizes  $\kappa(\mathbf{b})$ . Differentiating again with respect to  $\mathbf{b}$ , we have

$$\begin{aligned}\kappa'' &= -\sum_{i=1}^n \frac{\mathbf{z}_i \mathbf{z}_i^T}{\phi a_i v(\mu_i^b) [g'(\mu_i^b)]^2} + \mathbf{D}^{-1} + \mathbf{R} \\ &\approx \mathbf{Z}^T \mathbf{W} \mathbf{Z} + \mathbf{D}^{-1},\end{aligned}\tag{2.16}$$

where  $\mathbf{W}$  is the  $n \times n$  diagonal matrix with diagonal terms  $w_i = \{\phi a_i v(\mu_i^b) [g'(\mu_i^b)]^2\}^{-1}$  that are recognizable as the GLM iterated weights (Firth 1991, McCullagh & Nelder 1989). The remainder term  $\mathbf{R} = -\sum_{i=1}^n (y_i - \mu_i^b) \mathbf{z}_i \frac{\partial}{\partial \mathbf{b}} \left[ \frac{1}{\phi a_i v(\mu_i^b) g'(\mu_i^b)} \right]$  has expectation 0 and is thus, in probability as a function of  $n$ , of lower order than the two leading terms in the equation of  $\kappa''$ .  $\mathbf{R}$  equals  $\mathbf{0}$  for the canonical link functions, for which  $g'(\mu) = v^{-1}(\mu)$  (McCullagh & Nelder 1989). Combining (2.14)–(2.16) and ignoring  $\mathbf{R}$  leads to

$$ql(\boldsymbol{\alpha}, \boldsymbol{\theta}) \approx -\frac{1}{2} \log |\mathbf{I} + \mathbf{Z}^T \mathbf{W} \mathbf{Z} \mathbf{D}| - \frac{1}{2\phi} \sum_{i=1}^n d_i(y_i, \mu_i^{\tilde{\mathbf{b}}}) - \frac{1}{2} \tilde{\mathbf{b}}^T \mathbf{D}^{-1} \tilde{\mathbf{b}},\tag{2.17}$$

where  $\tilde{\mathbf{b}}$  is chosen to maximize the sum of the last two terms.

Assuming that the GLM iterative weights vary slowly (or not at all) as a function of the mean, the first term in this expression is ignored, and  $\boldsymbol{\alpha}$  is chosen to maximize the second. Thus  $(\hat{\boldsymbol{\alpha}}, \hat{\mathbf{b}}) = (\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}), \hat{\mathbf{b}}(\boldsymbol{\theta}))$ , where  $\hat{\mathbf{b}}(\boldsymbol{\theta}) = \tilde{\mathbf{b}}(\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}))$ , jointly maximize Green's (1987) PQL

$$-\frac{1}{2\phi} \sum_{i=1}^n d_i(y_i, \mu_i^{\mathbf{b}}) - \frac{1}{2} \mathbf{b}^T \mathbf{D}^{-1} \mathbf{b}.\tag{2.18}$$

Differentiation with respect to  $\boldsymbol{\alpha}$  and  $\mathbf{b}$  leads to the score equations for the mean parameters:

$$\begin{aligned}\sum_{i=1}^n \frac{(y_i - \mu_i^b) \mathbf{x}_i}{\phi a_i v(\mu_i^b) g'(\mu_i^b)} &= 0 \\ \sum_{i=1}^n \frac{(y_i - \mu_i^b) \mathbf{z}_i}{\phi a_i v(\mu_i^b) g'(\mu_i^b)} &= \mathbf{D}^{-1} \mathbf{b}.\end{aligned}$$

For the estimation of variance component  $\boldsymbol{\theta}$ , we substitute the maximized value of (2.18) into (2.17) and evaluate  $\mathbf{W}$  at  $(\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}), \hat{\mathbf{b}}(\boldsymbol{\theta}))$ , which generate an approximate profile quasi-likelihood function  $ql(\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}), \boldsymbol{\theta})$  for inference on  $\boldsymbol{\theta}$ . To make degrees-of-freedom adjustments that account for the fact that  $\hat{\boldsymbol{\alpha}}$  rather than  $\boldsymbol{\alpha}$  appears in the approximate profile quasi-likelihood function  $ql(\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}), \boldsymbol{\theta})$ , we modify  $ql(\hat{\boldsymbol{\alpha}}(\boldsymbol{\theta}), \boldsymbol{\theta})$  to the REML version (Patterson & Thompson 1971) in practice. By differentiating the modified profile quasi-likelihood with respect to the components of  $\boldsymbol{\theta}$ , we obtain the estimating equations for the variance parameters.

### 2.4.2 Bias correction in penalized quasi-likelihood

The approach proposed by Breslow & Clayton (1993) have been applied to a wide variety of generalized linear mixed models. Although the approximate procedure have been demonstrated to work reasonably well for discrete data problems with moderate to large cell frequencies, their performance is less satisfactory when the data are sparse. Breslow & Lin (1995) derived the general expressions for the asymptotic biases in approximate estimators of regression coefficients and variance component, for small values of the variance component, in generalized linear mixed models with canonical link function and a single source of extraneous variation. Their numerical studies of a series of matched pairs of binary outcomes showed that the first order estimators of the variance component are seriously biased, and they provided the easily computed correction factors which produce satisfactory estimators of small variance components. Their variance correction factors for a series of matched pairs of binomial observations rapidly approach one as the binomial denominators increase.

Let the data be in a series of  $m$  clusters of observations  $(y_{ij}, x_{ij})$ , where  $i$  identifies the cluster,  $j = 1, \dots, n_i$  identifies subjects within clusters and  $x_{ij}$  are  $p$ -vectors of explanatory variables associated with the univariate outcome  $y_{ij}$ . Given an unobserved

random effect  $b_i$ , the observations in the  $i$ -th cluster are assumed to have log conditional density

$$l_i(\alpha; b_i) = \sum_{j=1}^{n_i} \frac{a_{ij}}{\phi} \{y_{ij}\eta_{ij} - h(\eta_{ij})\} + c(y_{ij}; \phi), \quad (2.19)$$

where  $\eta_{ij} = x_{ij}^T \alpha + b_i$  denotes a linear predictor, the  $a_{ij}$  are prior weights, and  $\phi$  is a known scale parameter. This restriction to canonical link functions (McCullagh & Nelder, 1989) implies that the conditional means  $\mu_{ij}^{b_i} = E(y_{ij}|b_i) = g^{-1}(\eta_{ij})$  and variances  $\text{Var}(y_{ij}|b_i) = \phi a_{ij}^{-1} v(\mu_{ij}^{b_i})$  are related via  $g' = 1/h'' = 1/v$  for link and variance functions  $g$  and  $v$ , respectively. The  $b_i$  are assumed to be a random sample from a normal population with mean 0 and variance  $\theta$ . Thus the likelihood for the observed data is

$$L(\alpha, \theta) = \prod_{i=1}^m L_i(\alpha, \theta) = \prod_{i=1}^m (2\pi\theta)^{-\frac{1}{2}} \int e^{l_i(\alpha, b) - b^2/2\theta} db. \quad (2.20)$$

Denote  $(\hat{\alpha}, \hat{\theta})$  as the true maximum likelihood estimator. For approximations, we consider the derivatives  $l_i^{(k)} = \partial^k l_i / \partial b^k$ . Using Laplace method (e.g. Barndorff-Nielsen & Cox 1989), the likelihood function in (2.20) may be approximated by expanding the integrand in a Taylor series about its maximizing value  $\tilde{b}_i$ , where  $\tilde{b}_i = \tilde{b}_i(\alpha, \theta)$  solves  $\tilde{b}_i = \theta l_i^{(1)}(\alpha, \tilde{b}_i)$ . Setting  $\tilde{l}_i^{(k)} = l_i^{(k)}(\alpha, \tilde{b}_i)$ , a quartic expansion gives

$$\begin{aligned} L_i(\alpha, \theta) &\simeq (2\pi\theta)^{-\frac{1}{2}} \exp\left(\tilde{l}_i - \frac{\tilde{b}_i^2}{2\theta}\right) \int \exp\left\{\frac{1}{2}\left(\tilde{l}_i^{(2)} - \frac{1}{\theta}\right)(b - \tilde{b}_i)^2\right\} \\ &\quad \times \left\{1 + \frac{1}{6}\tilde{l}_i^{(3)}(b - \tilde{b}_i)^3 + \frac{1}{24}\tilde{l}_i^{(4)}(b - \tilde{b}_i)^4\right\} db \\ &= (1 - \theta\tilde{l}_i^{(2)})^{-\frac{1}{2}} \exp\left(\tilde{l}_i - \frac{\tilde{b}_i^2}{2\theta}\right) \left\{1 + \frac{\theta^2\tilde{l}_i^{(4)}}{8(1 - \theta\tilde{l}_i^{(2)})^2}\right\} \\ &\simeq (1 - \theta\tilde{l}_i^{(2)})^{-\frac{1}{2}} \exp\left\{\tilde{l}_i - \frac{\tilde{b}_i^2}{2\theta} + \frac{\theta^2\tilde{l}_i^{(4)}}{8(1 - \theta\tilde{l}_i^{(2)})^2}\right\}, \end{aligned}$$

where we evaluated the integral by taking expectations with respect to a normal variate having mean  $\tilde{b}_i$  and variance  $\theta/(1 - \theta\tilde{l}_i^{(2)})$ . We define the first order Laplace approxi-

mation to the log likelihood using only the leading terms of this expansion,

$$l_{L1}(\alpha, \theta) = \sum_{i=1}^m \left\{ -\frac{1}{2} \log(1 - \theta \tilde{l}_i^{(2)}) + \tilde{l}_i - \frac{\tilde{b}_i^2}{2\theta} \right\}. \quad (2.21)$$

The Laplace approximation estimator  $(\hat{\alpha}_{L1}, \hat{\theta}_{L1})$  are defined to be those that maximize  $l_{L1}$ .

Breslow & Clayton (1993), following Green (1987), termed the log conditional likelihoods (2.19) minus the penalty term  $\sum_i b_i^2/(2\theta)$  a log penalized quasi-likelihood in recognition of the fact that  $l_i$  requires specification only of the mean-variance relationship for the conditional distribution. Maximizing the penalized quasi-likelihood as a function of  $b = (b_1, \dots, b_m)^T$  for fixed  $(\alpha, \theta)$  leads to an objective function

$$l_p(\alpha, \theta) = \sum_{i=1}^m \left( \tilde{l}_i - \frac{\tilde{b}_i^2}{2\theta} \right)$$

that equals the sum of the last two terms in the first order Laplace approximation (2.21). The penalized quasi-likelihood estimator of the regression coefficients is defined to be the value  $\hat{\alpha}_P(\theta)$  that maximizes  $l_i(\alpha, \theta)$  for fixed  $\theta$ . The optimization may be programmed as a problem in iterated weighted least squares. Specifically, let  $\mathbf{Y}$  denote the  $N = \sum_i n_i$  dimensional ‘working vector’ whose components in lexicographic order are  $Y_{ij} = x_{ij}^T \alpha + \tilde{b}_i + (y_{ij} - \mu_{ij}^{\tilde{b}_i})/v_{ij}^{\tilde{b}_i}$ ; let  $\mathbf{V}$  denote the  $N \times N$  block diagonal covariance matrix whose  $n_i \times n_i$  dimensional diagonal submatrices  $V_i$  have terms  $\phi(a_{ij} v_{ij}^{\tilde{b}_i})^{-1} + \theta$  along their diagonals and off-diagonal elements  $\theta$ ; and let  $\mathbf{X}$  denote the  $N \times p$  design matrix with rows  $x_{ij}^T$ . Then, the Fisher scoring algorithm for solving the penalized quasi-likelihood equations

$$\frac{\partial l_p(\alpha, \theta)}{\partial \alpha} = \sum_{i=1}^m \sum_{j=1}^{n_i} \frac{a_{ij}}{\phi} (y_{ij} - \mu_{ij}^{\tilde{b}_i}) x_{ij} = 0 \quad (2.22)$$

for  $\alpha$  reduces to iterative solution of  $(\mathbf{X}^T \mathbf{V}^{-1} \mathbf{X})\alpha = \mathbf{X}^T \mathbf{V}^{-1} \mathbf{Y}$  (Green, 1987).

Estimation of  $\theta$  under penalized quasi-likelihood treats the working vector  $\mathbf{Y}$  as normally distributed with covariance matrix  $\mathbf{V}$  depending on  $\theta$ , except that the dependence of the terms  $v_{ij}^{\tilde{b}_i}$  on  $\theta$  through  $\tilde{b}_i$  is ignored when calculating derivatives (Breslow & Clayton, 1993). The principal advantage of this approach over the Laplace approximations is that it may be implemented using standard software for mixed model analysis. In this paper by Breslow & Lin (1995) the simpler maximum likelihood is used since they focused on asymptotic results rather than small sample properties while Breslow & Clayton (1993) used the restricted maximum likelihood normal theory approach. Thus, the penalized quasi-likelihood variance estimating equation is

$$\begin{aligned} \tilde{U}(\theta) &= \frac{1}{2} \left\{ (\mathbf{Y} - \mathbf{X}\alpha)^T \mathbf{V}^{-1} \frac{\partial \mathbf{V}}{\partial \theta} \mathbf{V}^{-1} (\mathbf{Y} - \mathbf{X}\alpha) - \text{tr} \left( \mathbf{V}^{-1} \frac{\partial \mathbf{V}}{\partial \theta} \right) \right\} \Big|_{\alpha = \hat{\alpha}_P(\theta)} \\ &= \frac{1}{2} \sum_{i=1}^m \left( \tilde{l}_i^{(1)2} + \frac{\tilde{l}_i^{(2)}}{1 - \theta \tilde{l}_i^{(2)}} \right) \Big|_{\alpha = \hat{\alpha}_P(\theta)} = 0. \end{aligned} \quad (2.23)$$

The penalized quasi-likelihood estimators  $(\hat{\alpha}_P, \hat{\theta}_P)$  simultaneously solve equations (2.22) and (2.23). While  $\hat{\alpha}_P(\theta)$  maximizes  $l_P(\alpha, \theta)$ , however,  $\hat{\theta}_P$  does not maximize  $l_P\{\hat{\alpha}_P(\theta), \theta\}$ .

Depending upon the distribution of the data and thus the link function in canonical generalized linear mixed models, the estimates of regression coefficients may be heavily influenced by the value assumed for the dispersion parameter. Accordingly, since some of the bias in an estimator of  $\alpha$  may arise from bias in the corresponding estimator of  $\theta$ , Breslow & Lin (1995) studied the bias in the estimator of  $\alpha$  for small fixed  $\theta$ , and then the bias in the estimator of  $\theta$ .

First, we expand the true log-likelihood and the approximation in Taylor series

about  $\theta = 0$ . Then, we have

$$\begin{aligned} l &= \sum_{i=1}^m l_{i0} + \theta \sum_{i=1}^m \left( \frac{l_{i0}^{(2)}}{2} + \frac{l_{i0}^{(1)2}}{2} \right) + \frac{\theta^2}{2} \sum_{i=1}^m \left( \frac{l_{i0}^{(2)2}}{2} + l_{i0}^{(1)2} l_{i0}^{(2)} + l_{i0}^{(1)} l_{i0}^{(3)} + \frac{l_{i0}^{(4)}}{4} \right) + o(\theta^2) \\ l_P &= l - \frac{\theta}{2} \sum_{i=1}^m l_{i0}^{(2)} - \frac{\theta^2}{4} \sum_{i=1}^m l_{i0}^{(2)2} - \frac{\theta^2}{2} \sum_{i=1}^m l_{i0}^{(1)} l_{i0}^{(3)} - \frac{\theta^2}{8} \sum_{i=1}^m l_{i0}^{(4)} + o(\theta^2), \end{aligned}$$

where we use the fact that

$$\left. \frac{\partial \tilde{b}_i}{\partial \theta} \right|_{\theta=0} = l_{i0}^{(1)}, \quad \left. \frac{\partial \tilde{b}_i}{\partial \theta^2} \right|_{\theta=0} = 2l_{i0}^{(1)} l_{i0}^{(2)}.$$

Then, the difference between  $\hat{\alpha}_P$  and  $\hat{\alpha}$  are studied by expanding

$$0 = \left. \frac{\partial l}{\partial \alpha} \right|_{\alpha=\hat{\alpha}} = \left. \frac{\partial l}{\partial \alpha} \right|_{\alpha=\hat{\alpha}_P} + \left. \frac{\partial^2 l}{\partial \alpha \alpha^T} \right|_{\alpha=\alpha^*} (\hat{\alpha} - \hat{\alpha}_P).$$

Consequently, we have

$$\hat{\alpha}_P = \hat{\alpha} + \frac{\theta}{2} (\mathbf{X}^T \mathbf{W}_0 \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u} + o(\theta), \quad (2.24)$$

where  $\mathbf{W}_0$  denotes the diagonal matrix with weight  $a_{ij} v_{ij}^0 / \phi$  on the diagonal  $\mathbf{u}$  is an  $N \times 1$  vector with components  $a_{ij} v(\mu_{ij}^0) v'(\mu_{ij}^0) / \phi$  and both  $\partial l / \partial \alpha$  and  $\partial l_P / \partial \alpha$  are evaluated at  $\alpha = \hat{\alpha}_P(\theta)$ . The corrected penalized quasi-likelihood estimate is obtained by subtracting the linear term in (2.24) from  $\hat{\alpha}_P$ .

The asymptotic biases in the estimator of  $\theta$  derived from the penalized quasi-likelihood were evaluated by equating expansions of the log profile likelihood  $l^\#(\theta) = \log L\{\hat{\alpha}(\theta), \theta\}$  to expansion of the penalized quasi-likelihood approximations. Then, we have

$$\frac{\hat{\theta}_P}{\hat{\theta}} = \left( \frac{\partial^2 l_P^\#}{\partial \theta^2} \right)^{-1} \left. \frac{\partial^2 l^\#}{\partial \theta^2} \right|_{\theta=0} \sim \frac{B - C}{C},$$

where

$$\begin{aligned}
B &= \sum_i l_{i0}^{(2)2}/2 - \mathbf{u}^T \mathbf{X} (\mathbf{X}^T \mathbf{W}_0 \mathbf{X})^{-1} \mathbf{X}^T \mathbf{u}/4, \\
C &= \sum_i l_{i0}^{(4)}/4, \\
D &= \sum_i l_{i0}^{(2)2}/2.
\end{aligned}$$

Lin & Brelsow (1996) also derived the biases correction in generalized linear mixed models with multiple components of dispersion.



# Chapter 3

## JOINT ANALYSIS FOR SURVIVAL TIME AND LONGITUDINAL CATEGORICAL MEASUREMENTS OF QUALITY OF LIFE IN HEAD AND NECK CANCER PATIENTS

### 3.1 Introduction

When choosing a treatment, it is well known that decision-making on treatment is frequently based on probability of survival. However, when there are multiple treatment modalities with similar survival rates, Quality of Life (QoL) factors are raised as important considerations for patients. In particular, oncology community has recognized that QoL and functional status are the major outcome variables in the evaluation of head and neck cancer treatment because of the potential impact on critical functions such as speech, swallowing, and breathing, as well as cosmesis and communication.

Many studies have been conducted for QoL without considering the survival time. For example, Terrell *et al.* (2004) investigated clinical predictors of QoL in a large intervention study of patients with head and neck cancer. Ringash, Bezjak, O'Sullivan

and Redelmeier (2004) studied QoL of particularly laryngeal cancer patients among those with head and neck cancer. Holloway *et al.* (2005) studied psychosocial effects in long-term head and neck cancer survivors. Fang *et al.* (2004) studied changes in QoL of head and neck cancer patients following postoperative radiotherapy. Most recently, Nibu *et al.* (2010) collected QoL data at scheduled clinic appointments of head and neck cancer patients and conducted a longitudinal QoL analysis. All these studies did not take the survival time into consideration. In order to completely understand the factors influencing both QoL and survival, it is important to study the QoL and survival simultaneously.

The Carolina Head and Neck Cancer Study (CHANCE) is a population based epidemiologic study conducted at 60 hospitals in 46 counties in North Carolina from 2002 through 2006 (Divaris *et al.* 2010). Patients were diagnosed with head and neck cancer (oral, pharynx, and larynx cancer) from 2002–2006. Their survival status was collected up to 2007 and QoL was evaluated over time for three years after diagnosis. QoL information was collected through questionnaires. Based on summary scores of the five domains of self-perceived quality of life including Physical Well-Being (PWB), Social/Family Well-Being (SWB), Emotional Well-Being (EWB), Functional Well-Being (FWB) and Head and Neck Cancer Specific symptoms (HNCS), patient’s QoL information was classified into satisfaction or dissatisfaction with life. Survival time is defined as the time to death from diagnosis. Demographic and life style characteristics, medical histories and clinical factors are also collected. It is of interest to elucidate the variables which are associated with both QoL satisfaction and survival time for patients with head and neck cancer. Additionally, the longitudinal QoL satisfaction outcomes and survival time are correlated within a patient, and this dependency should be taken into account in the analysis.

Among the existing approaches for longitudinal data and survival time, the selec-

tion model and the pattern mixture model have been widely used. The selection model estimates the distribution of survival time given longitudinal data. The selection model with continuous longitudinal data was studied by Tsiatis, Degruittola, and Wulfsohn (1995), Wulfsohn and Tsiatis (1997), Henderson, Diggle and Dobson (2000), Tsiatis and Davidian (2001), Song, Davidian and Tsiatis (2002), Tseng, Hsieh and Wang (2005) and Song and Wang (2007). The selection model with categorical longitudinal data was considered by Faucett, Schenker and Elashoff (1998), Huang *et al.* (2001), Xu and Zeger (2001a,b) and Larsen (2004). The pattern mixture model focuses on the trend of longitudinal outcomes conditional on survival time. The pattern mixture model with continuous longitudinal outcomes was studied by Wu and Carroll (1988), Wu and Bailey (1989), Hogan and Laird (1997), Ribaud, Thompson and Allen-Merish (2000) and more recently Ding and Wang (2008). Albert and Follmann (2000) proposed to model repeated count data subject to informative dropout and Albert, Follmann, Wang and Suh (2002) and Albert and Follmann (2007) studied binary longitudinal data with informative missingness. These methods cannot be applied directly to assess covariate effects on both outcomes. Simultaneous modeling of the longitudinal and survival data are needed for such purpose. Xu and Zeger (2001b) and Zeng and Cai (2005a) proposed simultaneous models of continuous longitudinal outcome and survival time. In their articles, heterogeneity caused by unobserved factors is represented using subject-specific random effects. Given random effects, survival time and the repeated measurements of longitudinal outcomes are assumed to follow a Cox proportional hazards model and a Gaussian distribution, respectively. Recently, Elashoff, Li and Ni (2007, 2008) proposed a more general joint model which incorporates a competing risks model for survival endpoint. Rizopoulos, Verbeke and Molenberghs (2008) assumed an accelerated failure time model and proposed to consider two separate sets of random effects for the continuous longitudinal process and survival time process, linking them

using a copula function. As an extension of this study, Rizopoulos, Verbeke, Lesaffre and Vanrenterghem (2008) considered longitudinal binary data with excess zeros and proposed a two-part shared parameter model framework. In the Bayesian perspective, Wang and Taylor (2001) and Brown and Ibrahim (2003) studied the simultaneous analysis of continuous longitudinal outcomes and survival time. Hu, Li and Li (2009) extended the existing Bayesian approach by considering the more general joint model of Elashfoff *et al.* (2008) with multiple types of failures in the failure time data.

Compared to the studies for continuous longitudinal data and survival time, relatively little work has been done in the joint modeling framework for categorical longitudinal data and survival time. However, the outcomes may not be continuous in some biomedical studies, for example, where the outcomes are disease symptom with categories of mild/moderate/severe, quality of life measurements with dissatisfied/satisfied, or dichotomized test results with categories of positive/negative. With these categorical longitudinal outcomes, the existing theory cannot be applied directly and the numerical algorithm needs to be modified. Therefore, in this paper, we investigate the simultaneous modeling of survival time and longitudinal categorical outcomes. Furthermore, hazards model for survival time is extended to allow multiple strata in our approach. Random effects are introduced into the proposed models to account for the dependence between survival time and longitudinal outcomes due to unobserved factors.

The outline of this paper is as follows. We begin by describing the details of the CHANCE study in Section 3.2. In Section 3.3, we then present a simultaneous modeling for longitudinal categorical outcomes and survival time, and describe the inference procedure. Asymptotic properties of the proposed estimators and the technical details of their proofs are given in Section 3.4 and Section 3.5, respectively. Numerical results from simulation studies are given in Section 3.6. The analysis of the CHANCE study

is provided in Section 3.7. In Section 3.8, we discuss some further consideration and generalization.

## **3.2 The CHANCE Study**

The Carolina Head and Neck Cancer Study (CHANCE) is the largest epidemiologic study of squamous cell carcinoma of the head and neck in the United States and the first to include a significant number of black patients. Patients who were diagnosed with head and neck cancer (oral, pharynx, and larynx cancer) from 2002 to 2006 were evaluated for Quality of Life (QoL) at maximum three times over follow-up at one to six months, one year and three years after diagnosis. At each evaluation, they were given questionnaires asking about their QoL satisfaction. Ending in December 2009, information on QoL has been obtained from 587 head and neck cancer patients. Based on the death information through 2007 available from the National Death Index (NDI), 91 patients died. It is of interest to study the effects of demographic and life style characteristics, medical histories, and clinical factors on patients' QoL and survival time. In particular, it is of interest to compare between African-Americans and Whites since it is known that African-Americans have a higher incidence of head and neck cancer and worse survival than Whites. Furthermore, because QoL outcomes are especially critical for physicians, head and neck cancer patients, and their caregivers, more research was needed on the experiences of survivors, especially among black patients. Given the paucity of data and studies on QoL among African-American head and neck cancer survivors, this study yields valuable new data.

To collect QoL information, the Functional Assessment of Cancer Therapy–Head and Neck Version 4 (FACT–H&N) series of questionnaires was used (Cella *et al.* 1993; Cella 1994; D'Antonio, Zimmerman, Cella and Long 1996; List *et al.* 1996). This QoL instrument was specifically designed for use of head and neck cancer patients and

consists of five primary QoL domains: Physical Well-Being (PWB), Social/Family Well-Being (SWB), Emotional Well-Being (EWB), Functional Well-Being (FWB) and Head and Neck Cancer Specific symptoms (HNCS). FACT-HN is the overall measurement of QoL incorporating all these domains. Each question has 5 scales of QoL measurement. Among them, the 3 high levels of “very much”, “quite a bit” and “somewhat” are categorized to “satisfied” and the other 2 low levels of “a little bit” and “not at all” belong to “dissatisfied”. We are interested in identifying the factors associated with both QoL and survival time. Longitudinal QoL outcomes are binary measurements with 1 (“satisfied”) and 0 (“dissatisfied”) on the five QoL domains, and survival time is the time to death from diagnosis. In this study, 33 among 587 patients were excluded due to missing data on household income, beer and QoL information resulting in 554 patients in the analysis. Eighty-five patients deceased by the end of 2007 and the censoring rate is 85%. Table 3.1 shows the descriptive statistics of predictors: demographics factors – race, household income, age at diagnosis, number of persons supported by household income; alcohol factor – the number of 12 oz. beers consumed per week; medical history factors – BMI and total number of medical conditions reported; treatment history factors – surgery, radiation therapy and chemotherapy; primary tumor data factors – primary tumor site and stage; time from diagnosis to each QOL survey. We analyze a QoL domain of the Head and Neck Cancer Specific symptoms (HNCS) and survival time, and Table 3.2 shows the descriptive statistics of outcome variables: longitudinal HNCS binary outcomes at three surveys, survival time from diagnosis and censorship indicator. The number of observations per patient ranges 1 to 3 with average of 1.93. We are interested in investigating factors which are associated with QoL and survival. In the next section, we formulate a general model and propose an inference procedure.

Table 3.1: Descriptive statistics of predictors in the CHANCE study

Categorical variables	Freq.	%
Total	554	100.00
Race		
– White	444	80.14
– African-American	110	19.86
Household income		
– level1: 0–10K	157	28.86
– level2: 20–30K	129	23.71
– level3: 40–50K	107	19.67
– level4: $\geq$ 60K	151	27.76
# of 12 oz. beers consumed per week		
– None	103	18.59
– less than 1	50	9.03
– 1 to 4	94	16.97
– 5 to 14	129	23.29
– 15 to 29	69	12.45
– 30 or more	109	19.68
Surgery		
– No	237	42.78
– Yes	317	57.22
Radiation therapy		
– No	131	23.65
– Yes	423	76.35
Chemotherapy		
– No	324	58.48
– Yes	230	41.52
Tumor site		
– Oral & Pharyngeal	346	62.45
– Laryngeal	208	37.55
Tumor stage		
– I	144	25.99
– II	93	16.79
– III	99	17.87
– IV	218	39.35

Continuous variables	n	mean	std.dev	min	median	max
Age at diagnosis	554	59.11	10.19	24.00	59.00	80.00
# of persons supported by household income	554	2.23	1.06	1.00	2.00	5.00
BMI	554	27.47	5.98	15.66	26.48	56.28
Total # of medical conditions reported	554	.92	1.10	.00	1.00	6.00
Time at 1st survey measurement (years)	209	.41	.45	.09	.28	3.55
Time at 2nd survey measurement (years)	500	1.85	.86	.44	1.81	3.91
Time at 3rd survey measurement (years)	353	3.49	.54	1.88	3.54	4.88

– Time at survey measurement is from diagnosis.

Table 3.2: Descriptive statistics of outcome variables in the CHANCE study

Longitudinal QoL binary outcomes	1st survey		2nd survey		3rd survey	
	Freq.	%	Freq.	%	Freq.	%
HNCS	209	100.00	500	100.00	352	100.00
– Dissatisfied (=0)	81	38.76	120	24.00	72	20.45
– Satisfied (=1)	128	61.24	380	76.00	280	79.55

Survival outcomes	n	mean	std.dev	min	median	max
min(Survival time, Censored time) (years)	554	3.07	1.04	.44	2.91	5.98
	Freq.	%				
Censorship	554	100.00				
– Alive	469	84.66				
– Death	85	15.34				

### 3.3 Models and Inference Procedure

#### 3.3.1 Model formulation and notation

Longitudinal measurements are considered as the realizations of a certain marker process at finite time points, and we use  $Y(t)$  to denote the value of such a marker process at time  $t$ . We let  $T$  be survival time, and suppose that the survival time  $T$  is possibly right censored and the right-censoring time is missing at random. Suppose a set of  $n$  subjects are followed over an interval  $[0, \tau]$ , where  $\tau$  is the study end time. Denote  $\mathbf{b}_i$ ,  $i = 1, \dots, n$ , as a vector of subject-specific random effects of dimension  $d_b$  and  $\mathbf{b}_i$ 's are mutually independent and identically distributed from a multivariate normal with mean zero and covariance matrix  $\Sigma_b$ .

Given the random effects  $\mathbf{b}_i$ , the observed covariates, and the observed outcome history till time  $t$ , we assume that the longitudinal outcome  $Y_i(t)$  at time  $t$  for subject



$i$  follows a distribution from the exponential family with density

$$\exp \left\{ \frac{y_i \eta_i(t) - B(\eta_i(t))}{A(D_i(t; \phi))} + C(y_i, D_i(t; \phi)) \right\} \quad (3.1)$$

with  $\mu_i(t) = E(Y_i(t)|\mathbf{b}_i) = B'(\eta_i(t))$  and  $v_i(t) = \text{Var}(Y_i(t)|\mathbf{b}_i) = B''(\eta_i(t))A(D_i(t; \phi))$ , satisfying

$$\eta_i(t) = g(\mu_i(t)) = \mathbf{X}_i(t)\boldsymbol{\beta} + \tilde{\mathbf{X}}_i(t)\mathbf{b}_i$$

and  $v_i(t) = v(\mu_i(t))A(D_i(t; \phi))$ , where  $g(\cdot)$  and  $v(\cdot)$  are known link and variance functions respectively, and  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  are the row vectors of the observed covariates for subject  $i$ , and  $\boldsymbol{\beta}$  is a column vector of coefficients for  $\mathbf{X}_i(t)$ . The random effect  $\mathbf{b}_i$  is allowed to differ for different individuals. Additionally,  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  can be completely different or share some components, and may include dummy variables for different strata.

Given the random effects  $\mathbf{b}_i$ , the observed covariates, and the observed survival history before time  $t$ , the conditional hazard rate function for the survival time  $T_i$  of subject  $i$  is assumed to follow a stratified multiplicative hazards model

$$\lambda_s(t) \exp\{\tilde{\mathbf{Z}}_i(t)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(t)\boldsymbol{\gamma}\}, \quad (3.2)$$

where  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  are the row vectors of the observed covariates and may share some components,  $\boldsymbol{\psi}$  is a vector of parameters of the coefficients for random effects,  $\lambda_s(t)$  is the  $s$ -th stratum baseline hazard rate function, and  $\boldsymbol{\gamma}$  is a column vector of coefficients for  $\mathbf{Z}_i(t)$ . Note that  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  do not include dummy variables for strata since baseline hazard rate is stratum-specific. Here, for any vectors  $\mathbf{a}_1$  and  $\mathbf{a}_2$  of the same dimension,  $\mathbf{a}_1 \circ \mathbf{a}_2$  denotes the component-wise product. In addition,  $\tilde{\mathbf{X}}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  have the same dimensions as  $\mathbf{b}_i$ 's.

Under models (3.1) and (3.2), the two outcomes  $Y(t)$  and  $T$  are independent conditional on the covariates and random effect. The parameter  $\boldsymbol{\psi}$  in model (3.2) characterizes the dependence between the longitudinal outcomes and the survival time due to latent random effect:  $\boldsymbol{\psi} = \mathbf{0}$  means that the dependence between the survival time and longitudinal responses are not due to these latent variables;  $\boldsymbol{\psi} \neq \mathbf{0}$  means that such dependence may be due to these latent variables. In other words,  $\boldsymbol{\psi} > \mathbf{0}$  implies that there may be some latent factors increasing both the longitudinal outcomes and the risk of survival endpoint simultaneously while  $\boldsymbol{\psi} < \mathbf{0}$  implies that some latent factors causing the increment of longitudinal outcomes may decrease the risk of survival endpoint.

We let  $n_i$  be the number of the observed longitudinal measurements for subject  $i$ , and assume that  $n_i$  and the observation times for longitudinal measurements are not informative about parameters of interest. That is, the distributions of  $n_i$  and the observation times for longitudinal measurements are independent of the parameters of interest in this joint model. The observed data from  $n$  subjects are  $(n_i, Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$ ,  $j=1, \dots, n_i$ ,  $i=1, \dots, n$ , and  $(V_i, \Delta_i, S_i, \{(\mathbf{Z}_i(t), \tilde{\mathbf{Z}}_i(t)) : t \leq V_i\})$ ,  $i=1, \dots, n$ , where for subject  $i$ ,  $(Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$  is the  $j$ -th observation of  $(Y_i(t), \mathbf{X}_i(t), \tilde{\mathbf{X}}_i(t))$ ,  $C_i$  is the right-censoring time,  $V_i = \min(T_i, C_i)$ ,  $S_i$  denotes the stratum, and  $\Delta_i = I(T_i \leq C_i)$ .

Our goal is to estimate and make inferences on the parameters  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \boldsymbol{\phi}^T, \text{Vec}(\boldsymbol{\Sigma}_b)^T, \boldsymbol{\psi}^T, \boldsymbol{\gamma}^T)^T$  and the baseline cumulative hazard functions with  $S$  strata,  $\boldsymbol{\Lambda}(t) = (\Lambda_1(t), \dots, \Lambda_S(t))^T$ , where  $\Lambda_s(t) = \int_0^t \lambda_s(u) du$ ,  $s = 1, \dots, S$ .  $\text{Vec}(\cdot)$  operator creates a column vector from a matrix by stacking the diagonal and upper-triangle elements of the matrix.

### 3.3.2 Inference procedure

For all  $n$  subjects, we write  $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_n^T)^T$ ,  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T$ ,  $\mathbf{V} = (V_1, \dots, V_n)^T$ , and  $\mathbf{b} = (\mathbf{b}_1^T, \dots, \mathbf{b}_n^T)^T$ . Then, the likelihood function of the complete data  $(\mathbf{Y}, \mathbf{V}, \mathbf{b})$  for

$(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  has the form,

$$\begin{aligned}
L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}) &= \prod_{s=1}^S \prod_{i=1}^n [f(\mathbf{Y}_i, V_i | \mathbf{b}_i) f(\mathbf{b}_i)]^{I(S_i=s)} = \prod_{i=1}^n f(\mathbf{Y}_i | \mathbf{b}_i) \left( \prod_{s=1}^S [f(V_i | \mathbf{b}_i)]^{I(S_i=s)} \right) f(\mathbf{b}_i) \\
&= \prod_{i=1}^n \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij} \boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i) - B(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \\
&\quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(V_i) \boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u) \boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right\},
\end{aligned}$$

and the full likelihood function of the observed data  $(\mathbf{Y}, \mathbf{V})$  for the parameter  $(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  is expressed as

$$L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = \int_{\mathbf{b}} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}) d\mathbf{b}. \quad (3.3)$$

The proposed estimation method is to calculate the maximum likelihood estimates for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda}(t))$  over a set in which  $\boldsymbol{\theta}$  is in a bounded set and  $\Lambda_s(t)$  of  $\boldsymbol{\Lambda}(t)$  belongs to a space consisting of all the increasing functions with  $\Lambda_s(0) = 0$ ,  $s = 1, \dots, S$ . We let each  $\Lambda_s(t)$  of  $\boldsymbol{\Lambda}(t)$ ,  $s = 1, \dots, S$ , be an increasing and right-continuous step function with jumps only at the observed failure times belonging to stratum  $s$ .

Denote  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}})$  as the maximum likelihood estimator for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda})$ . We let  $l_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}) = \log\{L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b})\}$  and  $l_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = \log\{L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V})\}$ , and denote  $U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  and  $U_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V})$  as the gradient vectors of the corresponding log-likelihood functions respectively. The EM-algorithm is used for calculating the maximum likelihood estimates. In the EM-algorithm,  $\mathbf{b}_i$  is considered as missing data for  $i = 1, \dots, n$ . Therefore, the M-step solves the conditional score equations from complete

data given observations, where the conditional expectation can be evaluated in E-step. The procedure involves iterating between the following two steps until convergence is achieved: at the  $k$ -th iteration,

(1) E-step Calculate the conditional expectations of some known functions of  $\mathbf{b}_i$ , needed in the next M-step, for subject  $i$  with  $S_i = s$  given observations and the current estimate  $(\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)})$ . To do this, denote  $q(\mathbf{b}_i)$  and  $E[q(\mathbf{b}_i)|\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}]$  as a known function and its conditional expectation, respectively. By some algebra,  $E[q(\mathbf{b}_i)|\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}]$  can be expressed in terms of a vector of new variables  $z_G$  following a multivariate Gaussian distribution with mean zero. The conditional expectation is calculated using the Gauss-Hermite Quadrature numerical approximation with 20 quadrature points.

(2) M-step After differentiating the conditional expectation of complete data log-likelihood function given observations and the current estimate  $(\boldsymbol{\theta}^{(k)}, \Lambda^{(k)})$ , the updated estimator  $(\boldsymbol{\theta}^{(k+1)}, \Lambda^{(k+1)})$  can be obtained as follows:  $(\boldsymbol{\beta}^{(k+1)}, \boldsymbol{\phi}^{(k+1)})$  solves the conditional expectation of complete data log-likelihood score equation using one-step Newton-Raphson iteration,

$$E[U_c(\boldsymbol{\beta}^{(k+1)}, \boldsymbol{\phi}^{(k+1)}|\boldsymbol{\theta}^{(k)}, \Lambda^{(k)})] = \mathbf{0},$$

where  $U_c(\boldsymbol{\beta}, \boldsymbol{\phi}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  is the first partial derivative of the complete data log-likelihood  $l_c(\boldsymbol{\theta}, \Lambda; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  with respect to  $(\boldsymbol{\beta}, \boldsymbol{\phi})$ ;

$$\Sigma_b^{(k+1)} = \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S E[\mathbf{b}_i \mathbf{b}_i^T | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}] I(S_i = s);$$

$(\boldsymbol{\psi}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)})$  solves the partial likelihood score equation from the full data using one-step Newton-Raphson iteration,

$$\begin{aligned}
& \sum_{i=1}^n \sum_{s=1}^S \Delta_i \left\{ \begin{pmatrix} E[(\tilde{\mathbf{Z}}_i^T(V_i) \circ \mathbf{b}_i) | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}] \\ \mathbf{Z}_i \end{pmatrix} \right. \\
& \quad \left. - \frac{\sum_{l: V_l \geq V_i} \begin{pmatrix} E[(\tilde{\mathbf{Z}}_l^T(V_i) \circ \mathbf{b}_l) \exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma} | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}\}] \\ E[\mathbf{Z}_l(V_i) \exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma} | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}\}] \end{pmatrix} I(S_l = s)}{\sum_{l: V_l \geq V_i} E[\exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma} | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}\}] I(S_l = s)} \right\} I(S_i = s) \\
& = \mathbf{0};
\end{aligned}$$

$\Lambda_s^{(k+1)}$  is obtained as an empirical function which has jumps only at the observed failure time,

$$\Lambda_s^{(k+1)}(t) = \sum_{i: V_i \leq t} \frac{\Delta_i I(S_i = s)}{\sum_{l: V_l \geq V_i} E[\exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi}^{(k+1)} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma}^{(k+1)} | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}\}] I(S_l = s)}.$$

The expressions of the conditional expectation and the conditional score equations calculated in the E- and M-steps for binary and Poisson longitudinal outcomes with survival time are given respectively in Sections 3.3.3.1 and 3.3.3.2.

The observed information matrix is adopted to obtain the variance estimate for  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}(t))$ . For the numerical calculation of the observed information matrix, we consider  $\Lambda_s\{V_i\}$ , the jump size of  $\Lambda_s(t)$  at  $V_i$  belonging to stratum  $s$  for which  $\Delta_i = 1$ , instead of  $\lambda_s(V_i)$ . That is,  $\boldsymbol{\Lambda}\{\cdot\} = (\boldsymbol{\Lambda}_1^T\{\cdot\}, \dots, \boldsymbol{\Lambda}_S^T\{\cdot\})^T$  with  $\boldsymbol{\Lambda}_s\{\cdot\} = (\Lambda\{T_{s1}\}, \dots, \Lambda\{T_{sm_s}\})^T$  for  $m_s$  failure times among  $n_s$  subjects ( $0 \leq m_s \leq n_s$ ) of the  $s$ -th stratum,  $s = 1, \dots, S$ . Then, by the Louis (1982) formula,

$$\begin{aligned}
I(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}) &= E_{b|Y,V}[B_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})|\mathbf{Y}, \mathbf{V}] \\
&- E_{b|Y,V}[U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})U_c^T(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})|\mathbf{Y}, \mathbf{V}] \\
&+ E_{b|Y,V}[U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V})]E_{b|Y,V}[U_c^T(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V})],
\end{aligned}$$

where  $B_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  is the negative of the second derivative matrix for the complete data log-likelihood  $l_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$ . The variance of  $\sqrt{n} \widehat{\boldsymbol{\theta}}$  is asymptotically equal to the corresponding sub-matrix of the inverse of the calculated observed information matrix. The variance of  $\widehat{\boldsymbol{\Lambda}}(t)$  is obtained using the estimated variances and covariances corresponding to  $\boldsymbol{\Lambda}\{\cdot\}$  from the inverse of the observed information matrix where  $T \leq t$  at the observed failures. In the EM-algorithm for variance estimation, we evaluate these conditional expectations only at the last iteration of the EM procedure for point estimation, where  $U_f$  is zero.

### 3.3.3 EM algorithm – examples

#### 3.3.3.1 Binary longitudinal data and survival time

(1) E-step : For binary longitudinal outcomes and survival time, we calculate the conditional expectation of  $q(\mathbf{b}_i)$  for subject  $i$  with  $S_i = s$  given the observations and the current estimate  $(\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)})$  for some known function  $q(\cdot)$ . The conditional expectation denoted by  $E[q(\mathbf{b}_i)|\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}]$  can be expressed as the following:

Given the current estimate  $(\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)})$ ,

$$E[q(\mathbf{b}_i)|\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}] = \frac{\int_{\mathbf{z}_G} q(R(\mathbf{z}_G))K(\mathbf{z}_G) \exp\{-\mathbf{z}_G^T \mathbf{z}_G\} d\mathbf{z}_G}{\int_{\mathbf{z}_G} K(\mathbf{z}_G) \exp\{-\mathbf{z}_G^T \mathbf{z}_G\} d\mathbf{z}_G}, \quad (3.4)$$

where

$$R(\mathbf{z}_G) = (\Sigma_b^{(k)})^{\frac{1}{2}} \left[ \sqrt{2} \mathbf{z}_G + (\Sigma_b^{(k)})^{\frac{1}{2}} \left( \sum_{j=1}^{n_i} y_{ij} \tilde{\mathbf{X}}_{ij}^T + \Delta_i(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(k)}) \right) \right], \quad (3.5)$$

$$K(\mathbf{z}_G) = \exp \left\{ - \sum_{j=1}^{n_i} \log \left( 1 + e^{\mathbf{X}_{ij} \boldsymbol{\beta}^{(k)} + \tilde{\mathbf{X}}_{ij} R(\mathbf{z}_G)} \right) - \int_0^{V_i} e^{\tilde{\mathbf{Z}}_i(u) (\boldsymbol{\psi}^{(k)} \circ R(\mathbf{z}_G)) + \mathbf{Z}_i(u) \boldsymbol{\gamma}^{(k)}} d\Lambda_s^{(k)}(u) \right\},$$

$(\Sigma_b^{(k)})^{\frac{1}{2}}$  is an unique non-negative square root of  $\Sigma_b^{(k)}$  (i.e.  $(\Sigma_b^{(k)})^{\frac{1}{2}} \times (\Sigma_b^{(k)})^{\frac{1}{2}} = \Sigma_b^{(k)}$ ),

and  $\mathbf{z}_G$  follows a multivariate Gaussian distribution with mean zero.

(2) M-step : Since the parameter  $\phi$  is set to 1 for logistic distribution, we estimate only  $\boldsymbol{\beta}$  in the longitudinal process.  $\boldsymbol{\beta}^{(k+1)}$  solves the conditional expectation of complete data log-likelihood score equation, using one-step Newton-Raphson iteration,

$$\begin{aligned} & \mathbb{E}[U_c(\boldsymbol{\beta}^{(k+1)} | \boldsymbol{\theta}^{(k)}, \Lambda^{(k)})] \\ &= \sum_{i=1}^n \sum_{j=1}^{n_i} \left( y_{ij} - \sum_{s=1}^S \mathbb{E} \left[ \frac{\exp\{\mathbf{X}_{ij} \boldsymbol{\beta}^{(k+1)} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i\}}{1 + \exp\{\mathbf{X}_{ij} \boldsymbol{\beta}^{(k+1)} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i\}} \middle| \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)} \right] I(S_i = s) \right) \mathbf{X}_{ij}^T = \mathbf{0}, \end{aligned}$$

where  $U_c(\boldsymbol{\beta}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  is the first partial derivative of the complete data log-likelihood  $l_c(\boldsymbol{\theta}, \Lambda; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  with respect to  $\boldsymbol{\beta}$ .  $\Sigma_b^{(k+1)}$ ,  $(\boldsymbol{\psi}^{(k+1)}, \boldsymbol{\gamma}^{(k+1)})$ , and  $\Lambda_s^{(k+1)}$  have the same expressions as in Section 3.3.2.

### 3.3.3.2 Poisson longitudinal data and survival time

(1) E-step : For Poisson longitudinal outcomes and survival time, given the current estimate  $(\boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)})$ , the conditional expectation denoted by  $E[q(\mathbf{b}_i) | \boldsymbol{\theta}^{(k)}, \Lambda_s^{(k)}]$  can be expressed as in (3.4) with  $R(\mathbf{z}_G)$  defined as in (3.5),

$$K(\mathbf{z}_G) = \exp \left\{ - \sum_{j=1}^{n_i} e^{\mathbf{X}_{ij} \boldsymbol{\beta}^{(k)} + \tilde{\mathbf{X}}_{ij} R(\mathbf{z}_G)} - \int_0^{V_i} e^{\tilde{\mathbf{Z}}_i(u) (\boldsymbol{\psi}^{(k)} \circ R(\mathbf{z}_G)) + \mathbf{Z}_i(u) \boldsymbol{\gamma}^{(k)}} d\Lambda_s^{(k)}(u) \right\},$$

and  $\mathbf{z}_G$  follows a multivariate Gaussian distribution with mean zero.

(2) M-step : Since the parameter  $\phi$  is set to 1 for Poisson distribution, we estimate only  $\beta$  in the longitudinal process.  $\beta^{(k+1)}$  solves the conditional expectation of complete data log-likelihood score equation, using one-step Newton-Raphson iteration,

$$\begin{aligned} & \mathbb{E}[U_c(\beta^{(k+1)}|\theta^{(k)}, \Lambda^{(k)})] \\ &= \sum_{i=1}^n \sum_{j=1}^{n_i} \left( y_{ij} - \sum_{s=1}^S \mathbb{E}[\exp\{\mathbf{X}_{ij}\beta^{(k+1)} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i\}|\theta^{(k)}, \Lambda_s^{(k)}] I(S_i=s) \right) \mathbf{X}_{ij}^T = \mathbf{0}, \end{aligned}$$

where  $U_c(\beta; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  is the first partial derivative of the complete data log-likelihood  $l_c(\theta, \Lambda; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  with respect to  $\beta$ .  $\Sigma_b^{(k+1)}$ ,  $(\psi^{(k+1)}, \gamma^{(k+1)})$ , and  $\Lambda_s^{(k+1)}$  have the same expressions as in Section 3.3.2.

### 3.4 Asymptotic Properties

To study the asymptotic properties of the proposed estimator  $(\hat{\theta}, \hat{\Lambda}(t))$  with  $\hat{\theta} = (\hat{\beta}^T, \hat{\phi}^T, \text{Vec}(\hat{\Sigma}_b)^T, \hat{\psi}^T, \hat{\gamma}^T)^T$  and  $\hat{\Lambda}(t) = (\hat{\Lambda}_1(t), \dots, \hat{\Lambda}_S(t))^T$ , we assume the following conditions below.

- (A1) The true parameter  $\theta_0 = (\beta_0^T, \phi_0^T, \text{Vec}(\Sigma_{b0})^T, \psi_0^T, \gamma_0^T)^T$  belongs to a known compact set  $\Theta$  which lies in the interior of the domain for  $\theta$ .
- (A2) The true baseline hazard rate function  $\lambda_0(t) = (\lambda_{10}(t), \dots, \lambda_{S0}(t))$  is bounded and positive in  $[0, \tau]$ , where  $\tau$  is the time of study end.
- (A3) For the censoring time  $C$ ,  $P(C \geq \tau | \mathbf{Z}, \tilde{\mathbf{Z}}, \mathbf{X}, \tilde{\mathbf{X}}) = P(C = \tau | \mathbf{Z}, \tilde{\mathbf{Z}}, \mathbf{X}, \tilde{\mathbf{X}}) > 0$ .
- (A4) For the number of observed longitudinal measurements per subject  $n_N$ ,  $P(n_N > d_b | \mathbf{X}, \tilde{\mathbf{X}}) > 0$  with probability one, and  $P(n_N \leq n_0) = 1$  for some integer  $n_0$ .
- (A5) Both  $\mathbf{X}^T \mathbf{X}$  and  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  are full rank with positive probability. Moreover, if there



exist constant vectors  $\mathbf{c}_1$  and  $\mathbf{c}_2$  such that, with positive probability, for any  $t$ ,  $\mathbf{Z}(t)\mathbf{c}_1 = \alpha_0(t)$  and  $\tilde{\mathbf{Z}}(t) \circ \mathbf{c}_2 = 0$  for a deterministic function  $\alpha_0(t)$ , then  $\mathbf{c}_1 = 0$ ,  $\mathbf{c}_2 = 0$ , and  $\alpha_0(t) = 0$ .

Assumption (A3) means that, by the end of the study, some proportion of the subjects will still be alive and censored at the study end time  $\tau$ , and thus the maximum right censoring time is equal to  $\tau$ . Assumption (A4) implies that some proportion of the subjects have at least  $d_b$  longitudinal observations, and there exists an integer  $n_0$  such that  $P(n_N \leq n_0) = 1$ . Consistency and asymptotic distribution of the proposed estimator are summarized in the following two theorems. We will present outlines of the proofs here. The detailed technical proofs are given in Section 3.5.

**Theorem 3.1.** *Under the assumptions (A1)~(A5), as  $n \rightarrow \infty$ , the maximum likelihood estimator  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}}(t))$  is consistent under the product norm of the Euclidean distance and the supreme norm on  $[0, \tau]$ . That is,  $\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} \|\hat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$ , a.s., where  $\|\hat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| = \sum_{s=1}^S |\hat{\Lambda}_s(t) - \Lambda_{s0}(t)|$ .*

Consistency of Theorem 3.1 can be proved by verifying the following three steps: First, we show that the maximum likelihood estimate  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}})$  exists. This can be achieved by showing that the jump size  $\Lambda_s\{V_i\}$ , with  $\Delta_i = 1$ , is finite. Second, we show that, with probability one,  $\hat{\Lambda}_s(\tau)$ ,  $s = 1, \dots, S$ , are bounded as  $n \rightarrow \infty$ . This can be proved by showing  $\log \hat{\Lambda}_s(\tau)$  is bounded. Third, given that the second step is true, by Helly's selection theorem (van der Vaart, 1998), we can choose a subsequence of  $\hat{\Lambda}_s(t)$  such that  $\hat{\Lambda}_s(t)$  weakly converges to some right-continuous monotone function  $\Lambda_s^*(t)$  with probability one. For any sub-sequence, we can find a further sub-sequence, still denoted as  $\hat{\boldsymbol{\theta}}$ , such that  $\hat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$ . Using empirical process formulation and relevant Donsker properties with parameter identifiability, we can show that  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$ ,  $s = 1, \dots, S$ . Based on these results, we can conclude that, with probability

one,  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\widehat{\boldsymbol{\Lambda}}_s(t)$  converges to  $\boldsymbol{\Lambda}_{s0}(t)$  in  $[0, \tau]$ ,  $s = 1, \dots, S$ . Moreover, since  $\boldsymbol{\Lambda}_{s0}(t)$  is right-continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$  almost surely.

**Theorem 3.2.** *Under the assumptions (A1)~(A5), as  $n \rightarrow \infty$ ,  $\sqrt{n}((\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T, (\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t))^T)^T$  weakly converges to a Gaussian random element in  $R^{d_\theta} \times \ell^\infty[0, \tau] \times \dots \times \ell^\infty[0, \tau]$ , and the estimator  $\widehat{\boldsymbol{\theta}}$  is asymptotically efficient, where  $d_\theta$  is the dimension of  $\boldsymbol{\theta}$  and  $\ell^\infty[0, \tau]$  is the normed space containing all the bounded functions in  $[0, \tau]$ .*

Once consistency is held, the conditions of Theorem 3.3.1 in van der Vaart and Wellner (1996), which implies the asymptotic normality of Theorem 3.2, are verified via the tools of empirical processes. These conditions are restated in Theorem 4 of Parner (1998). The smooth conditions in Theorem 4 of Parner (1998) can be verified using the regularity of the log-likelihood function in terms of model parameters and the Donsker properties of the score operators. In particular, in the invertibility condition of the information operator in Theorem 4 of Parner (1998), the verification of the one-to-one property of the information operator is specific to our proposed models and requires non-trivial work. Therefore, by Theorem 3.3.1 of van der Vaart and Wellner (1996),  $\sqrt{n}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\boldsymbol{\Lambda}}_s - \boldsymbol{\Lambda}_{s0})$  weakly converges to a Gaussian process, and by Proposition 3.3.1 in Bickel *et al.* (1993),  $\widehat{\boldsymbol{\theta}}$  is an efficient estimator for  $\boldsymbol{\theta}_0$ .

### 3.5 Technical Details – Proofs for Asymptotic Properties

In this section, we present the detailed technical proofs for the asymptotic properties of the proposed estimator  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}(t))$  with  $\widehat{\boldsymbol{\theta}} = (\widehat{\boldsymbol{\beta}}^T, \widehat{\boldsymbol{\phi}}^T, \text{Vec}(\widehat{\boldsymbol{\Sigma}}_b)^T, \widehat{\boldsymbol{\psi}}^T, \widehat{\boldsymbol{\gamma}}^T)^T$  and  $\widehat{\boldsymbol{\Lambda}}(t) = (\widehat{\boldsymbol{\Lambda}}_1(t), \dots, \widehat{\boldsymbol{\Lambda}}_S(t))^T$ . Meanwhile, the supplementary proofs needed to prove the asymptotic properties are provided in Section 3.5.3. From the full likelihood function

of observed data  $(\mathbf{Y}, \mathbf{V})$  for  $(\boldsymbol{\theta}, \Lambda)$ ,

$$\begin{aligned}
L_f(\boldsymbol{\theta}, \Lambda; \mathbf{Y}, \mathbf{V}) &= \int_{\mathbf{b}} L_c(\boldsymbol{\theta}, \Lambda; \mathbf{Y}, \mathbf{V}, \mathbf{b}) d\mathbf{b} \\
&= \prod_{i=1}^n \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i) - B(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \\
&\quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right\} d\mathbf{b},
\end{aligned}$$

we have the observed log-likelihood function

$$\begin{aligned}
\sum_{i=1}^n \log &\left[ \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i) - B(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \\
&\times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right\} d\mathbf{b} \Big].
\end{aligned}$$

Then, we obtain the following modified object function by replacing  $\lambda_s(V_i)$  with  $\Lambda_s\{V_i\}$  in the above expression where  $\Lambda_s\{V_i\}$  is the jump size of  $\Lambda_s(t)$  at the observed time  $V_i$  with  $\Delta_i = 1$ ,

$$\begin{aligned}
l_n(\boldsymbol{\theta}, \Lambda) &= \sum_{i=1}^n \log \left[ \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i) - B(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \\
&\quad \times \left( \prod_{s=1}^S \left[ \Lambda_s\{V_i\}^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right)
\end{aligned}$$

$$\times (2\pi)^{-d_b/2} |\Sigma_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}_i^T \Sigma_b^{-1} \mathbf{b}_i \right\} d\mathbf{b} \Big], \quad (3.6)$$

and  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  maximizes  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  over the space  $\{(\boldsymbol{\theta}, \boldsymbol{\Lambda}) : \boldsymbol{\theta} \in \Theta, \boldsymbol{\Lambda} \in \mathbb{W}_n \times \mathbb{W}_n \cdots \times \mathbb{W}_n\}$ , where  $\mathbb{W}_n$  consists of all the right-continuous step functions only; that is,  $\boldsymbol{\Lambda} = (\Lambda_1, \dots, \Lambda_S)^T, s = 1, \dots, S, \Lambda_s \in \mathbb{W}_n$ . For the proofs of both Theorem 3.1 and Theorem 3.2, the modified object function is used in the place of the observed log-likelihood function.

### 3.5.1 Proof of consistency

Consistency can be proved by verifying the following three steps: First, we show the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  exists. Second, we show that, with probability one,  $\widehat{\Lambda}_s(\tau), s = 1, \dots, S$ , are bounded as  $n \rightarrow \infty$ . Third, if the second step is true, by Helly's selection theorem (p9 of van der Vaart, 1998), we can choose a subsequence of  $\widehat{\Lambda}_s$  such that  $\widehat{\Lambda}_s$  weakly converges to some right-continuous monotone function  $\Lambda_s^*$  with probability one; that is, the measure given by  $\mu_s([0, t]) = \widehat{\Lambda}_s(t)$  for  $t \in [0, \tau]$  weakly converges to the measure given by  $\mu_s^*([0, t]) = \Lambda_s^*(t)$ . By choosing a sub-sequence, we can further assume  $\widehat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$ . Thus, in this third step, we show  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \boldsymbol{\theta}_{s0}, s = 1, \dots, S$ .

Once the three steps are completed, we can conclude that, with probability one,  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\widehat{\boldsymbol{\Lambda}}_s$  converges to  $\boldsymbol{\Lambda}_{s0}$  in  $[0, \tau], s = 1, \dots, S$ . However, since  $\boldsymbol{\Lambda}_{s0}$  is continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$  almost surely. Then, the proof of Theorem 3.1 will be done.

In the first step, we will show the existence of the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$ . Since  $\boldsymbol{\theta}$  belongs to a compact set  $\Theta$  by the assumption (A1), it is sufficient to show that  $\Lambda_s\{V_i\}$ , the jump size of  $\Lambda_s$  at  $V_i$  for which  $\Delta_i = 1$ , is finite. Since, for each

subject  $i$  with  $\Delta_i = 1$ ,

$$\begin{aligned} \Lambda_s\{V_i\} \exp \left\{ \int_0^{V_i} \exp \left\{ \tilde{\mathbf{Z}}_i(t)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}_i(t)\boldsymbol{\gamma} \right\} d\Lambda_s(t) \right\} \\ \leq \exp \left\{ -2(\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}) \right\} (\Lambda_s\{V_i\})^{-1}, \end{aligned}$$

we have that, from (3.6),

$$\begin{aligned} l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \leq \sum_{i=1}^n \log \int_{\mathbf{b}} \left[ \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b})}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \\ \times \left( \prod_{s=1}^S \left[ (\Lambda_s\{V_i\})^{-\Delta_i} \exp \left\{ -\Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right\} \right]^{I(S_i=s)} \right) \\ \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b} \right\} \Big] d\mathbf{b}. \end{aligned}$$

Thus, if  $\Lambda_s\{V_i\} \rightarrow \infty$  for some  $i$  with  $\Delta_i = 1$ , then  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \rightarrow -\infty$ , which is contradictory to that  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  is bounded. Therefore, we conclude that  $\Lambda_s\{\cdot\}$ , the jump size of  $\Lambda_s$  for stratum  $s$ , must be finite. By the conclusion and the assumption (A1), the maximum likelihood estimate  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}})$  exists.

In the second step, we will show that  $\hat{\Lambda}_s(\tau)$  is bounded as  $n$  goes to infinity with probability one. We define  $\hat{\zeta}_s = \log \hat{\Lambda}_s(\tau)$  and rescale  $\hat{\Lambda}_s$  by the factor  $e^{\hat{\zeta}_s}$ . Then, we let  $\tilde{\Lambda}_s$  denote the rescaled function; that is,  $\tilde{\Lambda}_s(t) = \hat{\Lambda}_s(t)/\hat{\Lambda}_s(\tau) = \hat{\Lambda}_s(t)e^{-\hat{\zeta}_s}$ . thus,  $\tilde{\Lambda}_s(\tau) = 1$ . To prove this second step, it is sufficient to show  $\hat{\zeta}_s$  is bounded. After some algebra in (3.6), we obtain that, for any  $\boldsymbol{\Lambda} \in \mathbb{W} \times \mathbb{W} \cdots \times \mathbb{W}$ ,

$$\begin{aligned} n^{-1}l_n(\hat{\boldsymbol{\theta}}, \boldsymbol{\Lambda}) &= \frac{1}{2} \sum_{i=1}^n \left[ \sum_{j=1}^{n_i} \left( \frac{Y_{ij} \mathbf{X}_{ij} \hat{\boldsymbol{\beta}}}{A(D_i(t_j; \hat{\phi}))} + C(Y_{ij}; D_i(t_j; \hat{\phi})) \right) + \sum_{s=1}^S \Delta_i (\mathbf{Z}_i(V_i) \hat{\boldsymbol{\gamma}}) I(S_i=s) \right. \\ &\quad - \frac{1}{2} \log \left\{ (2\pi)^{d_b} |\hat{\boldsymbol{\Sigma}}_b| \right\} + \frac{1}{2} \mathbf{M}_i^T \mathbf{M}_i - \frac{1}{2} \log |\hat{\boldsymbol{\Sigma}}_b| + \sum_{s=1}^S \Delta_i I(S_i=s) \log \Lambda_s\{V_i\} \\ &\quad \left. + \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\hat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \hat{\phi}))} \right\} \right] d\mathbf{b}_0 \right] \end{aligned}$$

$$\left. - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\Lambda_s(t) \right\} \Big] d\mathbf{b}_0 \Big],$$

where

$$\begin{aligned} \mathbf{M}_i &= \widehat{\boldsymbol{\Sigma}}_b^{1/2} \left( \sum_{j=1}^{n_i} \frac{Y_{ij} \widetilde{\mathbf{X}}_{ij}}{A(D_i(t_j; \widehat{\phi}))} + \sum_{s=1}^S I(S_i = s) \Delta_i(\widetilde{\mathbf{Z}}_i(V_i) \circ \widehat{\boldsymbol{\psi}}^T) \right)^T, \\ \mathbf{b}_0 &= \boldsymbol{\Sigma}_b^{-1/2} \mathbf{b} - \mathbf{M}_i, \end{aligned}$$

and

$$Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}}) = (\widetilde{\mathbf{Z}}_i(t) \circ \widehat{\boldsymbol{\psi}}^T) \widehat{\boldsymbol{\Sigma}}_b^{1/2} \mathbf{b}_0 + \mathbf{Z}_i(t) \widehat{\boldsymbol{\gamma}} + (\widetilde{\mathbf{Z}}_i(t) \circ \widehat{\boldsymbol{\psi}}^T) \widehat{\boldsymbol{\Sigma}}_b^{1/2} \mathbf{M}_i.$$

Thus, since  $0 \leq n^{-1} l_n(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}) - n^{-1} l_n(\widehat{\boldsymbol{\theta}}, \widetilde{\boldsymbol{\Lambda}})$  where  $\widehat{\boldsymbol{\Lambda}} = e^{\widehat{\boldsymbol{\xi}}} \circ \widetilde{\boldsymbol{\Lambda}}$ , it follows that

$$\begin{aligned} 0 &\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \left( \log e^{\widehat{\zeta}_s} \widetilde{\Lambda}_s - \log \widetilde{\Lambda}_s \right) \\ &\quad + \frac{1}{n} \sum_{i=1}^n \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right. \right. \\ &\quad \left. \left. - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0 \\ &\quad - \frac{1}{n} \sum_{i=1}^n \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right. \right. \\ &\quad \left. \left. - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0. \quad (3.7) \end{aligned}$$

According to the assumption (A2), there exist some positive constants  $C1$ ,  $C2$  and  $C3$  such that  $|Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})| \leq C1 \|\mathbf{b}_0\| + C2 \|\mathbf{Y}_i\| + C3$ . By denoting  $\mathbf{b}_0$  as a vector of variables following a standard multivariate normal distribution, from concavity of the logarithm function, in the third term of (3.7),

$$\begin{aligned} &\log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0 \\ &= (2\pi)^{d_b/2} \log E_{\mathbf{b}_0} \left[ \exp \left\{ -\sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \right] \\ &\geq (2\pi)^{d_b/2} \log E_{\mathbf{b}_0} \left[ \exp \left\{ -\sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - e^{C1 \|\mathbf{b}_0\| + C2 \|\mathbf{Y}_i\| + C3} \right\} \right] \end{aligned}$$

$$\begin{aligned}
&\geq (2\pi)^{d_b/2} \mathbb{E}_{\mathbf{b}_0} \left[ - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - e^{C_1 \|\mathbf{b}_0\| + C_2 \|\mathbf{Y}_i\| + C_3} \right] \\
&= - e^{C_2 \|\mathbf{Y}_i\| + C_4} - C_5,
\end{aligned}$$

where  $C_4$  and  $C_5$  are positive constants. Then, since it is easily verified that  $\mathbb{E}_{\mathbf{b}_0} \left[ \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} + e^{C_1 \|\mathbf{b}_0\| + C_2 \|\mathbf{Y}_i\| + C_3} \right] < \infty$ , by the strong law of large numbers and the assumption (A4), the third term of (3.7)

$$\begin{aligned}
&-\frac{1}{n} \sum_{i=1}^n \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\theta})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0 \\
&\leq \frac{1}{n} \sum_{i=1}^n (e^{C_2 \|\mathbf{Y}_i\| + C_4} + C_5) \triangleq C_6
\end{aligned}$$

can be bounded by some constant  $C_6$  from above. Then (3.7) becomes

$$\begin{aligned}
0 &\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
&\quad + \frac{1}{n} \sum_{i=1}^n \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right. \right. \\
&\quad \left. \left. - \sum_{s=1}^S e^{\widehat{\zeta}_s} \int_0^{V_i} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\theta})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0 + C_6 \\
&\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
&\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right. \right. \\
&\quad \left. \left. - \sum_{s=1}^S e^{\widehat{\zeta}_s} \int_0^{\tau} e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\theta})} d\widetilde{\Lambda}_s(t) \right\} \right] d\mathbf{b}_0 \\
&\quad + \frac{1}{n} \sum_{i=1}^n I(V_i \neq \tau) \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right\} \right] d\mathbf{b}_0 + C_6 \\
&\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
&\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right\} \right] d\mathbf{b}_0
\end{aligned}$$

$$\left. - \sum_{s=1}^S e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_0 + C_7, \quad (3.8)$$

where  $C_7$  is a constant. On the other hand, since, for any  $\Gamma \geq 0$  and  $x > 0$ ,  $\Gamma \log(1 + x/\Gamma) \leq \Gamma x/\Gamma = x$ , we have that  $e^{-x} \leq (1 + x/\Gamma)^{-\Gamma}$ . Therefore, in the second term of (3.8),

$$\begin{aligned} & \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \\ & \leq \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right\} \times \left\{ 1 + \frac{\sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t)}{\Gamma} \right\}^{-\Gamma} \\ & \leq \Gamma^\Gamma \times \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} \right\} \times \left\{ \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\}^{-\Gamma} \\ & = \Gamma^\Gamma \times \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \Gamma \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \right\} \\ & \quad \times \left\{ \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\}^{-\Gamma}. \end{aligned} \quad (3.9)$$

Since  $Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}}) \geq -C_1 \|\mathbf{b}_0\| - C_2 \|\mathbf{Y}_i\| - C_3$ ,

$$\begin{aligned} \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) & \geq \int_0^\tau e^{-C_1 \|\mathbf{b}_0\| - C_2 \|\mathbf{Y}_i\| - C_3} d\widetilde{\Lambda}_s(t) \\ & = e^{-C_1 \|\mathbf{b}_0\| - C_2 \|\mathbf{Y}_i\| - C_3} \times \{\widetilde{\Lambda}_s(\tau) - \widetilde{\Lambda}_s(0)\} \\ & = e^{-C_1 \|\mathbf{b}_0\| - C_2 \|\mathbf{Y}_i\| - C_3}. \end{aligned}$$

Thus, in (3.9),  $\left\{ \int_0^\tau e^{Q_{1i}(t, \mathbf{b}_0, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\}^{-\Gamma} \leq e^{C_1 \Gamma \|\mathbf{b}_0\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma}$ ,

and

$$(3.9) \leq \int_{\mathbf{b}_0} \left[ \Gamma^\Gamma \times \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} - \Gamma \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} + C_1 \Gamma \|\mathbf{b}_0\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma \right\} \right] d\mathbf{b}_0.$$



Therefore, (3.8) gives that

$$\begin{aligned}
0 &\leq C_7 + \frac{1}{n} \sum_{i=1}^n \Delta_i \left( \sum_{s=1}^S \widehat{\zeta}_s \right) + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \left\{ \Gamma^\Gamma \times \exp \left\{ -\Gamma \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right\} \right. \\
&\quad \times \left. \int_{\mathbf{b}_0} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_0^T \mathbf{b}_0 - \sum_{j=1}^{n_i} \frac{B(\widehat{\beta}; \mathbf{b}_0)}{A(D_i(t_j; \widehat{\phi}))} + C_1 \Gamma \|\mathbf{b}_0\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma \right\} \right] d\mathbf{b}_0 \right\} \\
&= C_7 + \frac{1}{n} \sum_{i=1}^n \Delta_i \left( \sum_{s=1}^S \widehat{\zeta}_s \right) - \frac{\Gamma}{n} \sum_{i=1}^n I(V_i = \tau) \left( \sum_{s=1}^S \widehat{\zeta}_s \right) + C_8(\Gamma), \tag{3.10}
\end{aligned}$$

where  $C_8(\Gamma)$  is a deterministic function of  $\Gamma$ . For the  $s$ -th stratum, (3.10) is that

$$0 \leq C_7 + \sum_{i=1}^n \Delta_i I(S_i = s) \widehat{\zeta}_s - \frac{\Gamma}{n} \sum_{i=1}^n I(V_i = \tau) I(S_i = s) \widehat{\zeta}_s + C_8(\Gamma).$$

By the strong law of large numbers,  $\sum_{i=1}^n I(V_i = \tau) I(S_i = s) / n \rightarrow P(V_i = \tau, S_i = s) > 0$ .

Then, we can choose  $\Gamma$  large enough such that  $\sum_{i=1}^n \Delta_i I(S_i = s) / n \leq (\Gamma / 2n) \sum_{i=1}^n I(V_i = \tau) I(S_i = s)$ . Thus, we obtain that

$$0 \leq C_7 + C_8(\Gamma) - \frac{\Gamma}{2n} \sum_{i=1}^n I(V_i = \tau) I(S_i = s) \widehat{\zeta}_s.$$

In other words,

$$\widehat{\zeta}_s \leq \frac{(C_7 + C_8(\Gamma))2n}{\Gamma \sum_{i=1}^n I(V_i = \tau) I(S_i = s)} \rightarrow \frac{(C_7 + C_8(\Gamma))2}{\Gamma P(V_i = \tau, S_i = s)}.$$

If we denote  $B_{s0} = \exp \{2(C_7 + C_8(\Gamma)) / (\Gamma P(V_i = \tau, S_i = s))\}$ , we conclude that  $\widehat{\Lambda}_s(\tau) \leq B_{s0}$ ,  $s = 1, \dots, S$ . Note that the above arguments hold for every sample in the probability space except a set with zero probability. Therefore, we have shown that, with probability one,  $\widehat{\Lambda}_s(\tau)$  is bounded for any sample size  $n$ .

In the third step, the goal of this step is to show that, if  $\widehat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$  and  $\widehat{\Lambda}_s$  weakly converges to  $\Lambda^*$  with probability one, then  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \boldsymbol{\theta}_{s0}$ ,  $s = 1, \dots, S$ . We set some preliminaries as the followings: For convenience, we omit the index  $i$  for

subject and use  $\mathbf{O}$  to abbreviate the observed statistics  $(\mathbf{Y}, \mathbf{X}, \tilde{\mathbf{X}}, V, \Delta, n_N, s)$  and  $\{\mathbf{Z}(t), \tilde{\mathbf{Z}}(t), 0 \leq t \leq V\}$  for a subject. By dropping  $(\lambda_s(V))^\Delta$  from the complete data likelihood function, we define that

$$\begin{aligned} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) &= \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta} + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} + C(Y_j; D(t_j; \phi)) \right] \right\} \\ &\quad \times \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}(V)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}(V)\boldsymbol{\gamma}] \right. \\ &\quad \left. - \int_0^V \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma} \} d\Lambda_s(t) \right\} \\ &\quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b} \right\}, \end{aligned}$$

$$\text{and} \quad Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) = \frac{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \exp \{ \tilde{\mathbf{Z}}(v)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}(v)\boldsymbol{\gamma} \} d\mathbf{b}}{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b}}.$$

Furthermore, for any measurable function  $f(\mathbf{O})$ , we use operator notation to define  $\mathbf{P}_n f = n^{-1} \sum_{i=1}^n f(\mathbf{O}_i)$  and  $\mathbf{P} f = \int f d\mathbf{P} = \mathbb{E}[f(\mathbf{O})]$ . Thus,  $\mathbf{P}_n f$  is the empirical measure from  $n$  i.i.d observations and  $\sqrt{n}(\mathbf{P}_n - \mathbf{P})$  is the empirical process based on these observations. We also define a class  $\mathcal{F} = \{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau], \boldsymbol{\theta} \in \Theta, \Lambda_s \in \mathbb{W}, \Lambda_s(0) = 0, \Lambda_s(\tau) \leq B_{s0}\}$ , where  $B_{s0}$  is the constant given in the second step and  $\mathbb{W}$  contains all nondecreasing functions in  $[0, \tau]$ . According to the result proved in Section 3.5.3.1,  $\mathcal{F}$  is P-Donsker.

Let  $m_s$  denote the number of subjects in stratum  $s$ ; i.e.  $n = \sum_{s=1}^S m_s$ .  $V_s$  and  $\Delta_s$  denote the observed time and censoring indicator for a subject belonging to stratum  $s$ , respectively. Thus,  $V_{sk}$  and  $\Delta_{sk}$  are the  $k$ -th subject observed time and censoring indicator in stratum  $s$ .

Now we start the proof of the third step. Since  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}})$  maximizes the function  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$ , where  $\boldsymbol{\Lambda} = (\Lambda_1, \dots, \Lambda_S)^T$  and  $\Lambda_s, s = 1, \dots, S$ , are any step functions with

jumps only at  $V_i$  belonging to stratum  $s$  for which  $\Delta_i = 1$ , we differentiate  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  with respect to  $\Lambda_s\{V_{sk}\}$  and obtain the following equation, satisfied by  $\widehat{\Lambda}_s$ ,

$$\widehat{\Lambda}_s\{V_{sk}\} = \frac{\Delta_{sk}}{m_s \mathbf{P}_{m_s} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) \right\} \Big|_{v=V_{sk}}}$$

Imitating the above equation, we also can construct another function, denoted by  $\bar{\Lambda} = (\bar{\Lambda}_1, \dots, \bar{\Lambda}_S)^T$  such that  $\bar{\Lambda}_s$ ,  $s = 1, \dots, S$ , are also step functions with jumps only at the the observed  $V_{sk}$  and the jump size  $\bar{\Lambda}_s\{V_{sk}\}$  is given by

$$\bar{\Lambda}_s\{V_{sk}\} = \frac{\Delta_{sk}}{m_s \mathbf{P}_{m_s} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_{sk}}}.$$

Equivalently,

$$\bar{\Lambda}_s(t) = \frac{1}{m_s} \sum_{k=1}^{m_s} \frac{I(V_{sk} \leq t) \Delta_{sk}}{\mathbf{P}_{m_s} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_{sk}}}.$$

Then, we claim  $\bar{\Lambda}_s(t)$  uniformly converges to  $\Lambda_{s0}(t)$  in  $[0, \tau]$ . To prove the claim, note

$$\begin{aligned} & \sup_{t \in [0, \tau]} \left| \bar{\Lambda}_s(t) - \mathbf{E} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_s}} \right] \right| \\ &= \sup_{t \in [0, \tau]} \left| \frac{1}{m_s} \sum_{k=1}^{m_s} \frac{I(V_{sk} \leq t) \Delta_{sk}}{\mathbf{P}_{m_s} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_{sk}}} \right. \\ & \quad \left. - \mathbf{P}_{m_s} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_s}} \right] \right. \\ & \quad \left. + \mathbf{P}_{m_s} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_s}} \right] \right. \\ & \quad \left. - \mathbf{P} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} \Big|_{v=V_s}} \right] \right| \end{aligned}$$

$$\begin{aligned}
&\leq \sup_{t \in [0, \tau]} \left| \frac{1}{m_s} \sum_{k=1}^{m_s} I(V_{sk} \leq t) \Delta_{sk} \left[ \frac{1}{\mathbf{P}_{m_s} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right. \right. \\
&\quad \left. \left. - \frac{1}{\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right] \right|_{v=V_{sk}} \\
&\quad + \sup_{t \in [0, \tau]} \left| (\mathbf{P}_{m_s} - \mathbf{P}) \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right] \right|_{v=V_s} \\
&\leq \sup_{t \in [0, \tau]} \left| \frac{1}{\mathbf{P}_{m_s} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} - \frac{1}{\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right| \\
&\quad + \sup_{t \in [0, \tau]} \left| (\mathbf{P}_{m_s} - \mathbf{P}) \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right] \right|. \tag{3.11}
\end{aligned}$$

In (3.11), the right hand side converges to 0 because the first and second terms on the right hand side converges to 0 in the following: First, according to Section 3.5.3.1,  $\{Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) : v \in [0, \tau]\}$  is a bounded and Glivenko-Cantelli class.  $\{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) : v \in [0, \tau]\}$  is also a Glivenko-Cantelli class because  $\{I(V_s \geq v) : v \in [0, \tau]\}$  is a Glivenko-Cantelli class and the functional  $(f, g) \rightarrow fg$  for any bounded two functions  $f$  and  $g$  is Lipschitz continuous. Then, we obtain that  $\sup_{t \in [0, \tau]} \left| \mathbf{P}_{m_s} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} - \mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \right|$  converges to 0. Besides, from Section 3.5.3.1,  $\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} > \mathbf{P} \{I(V_s \geq v) \exp\{-C_9 - C_{10} \|\mathbf{Y}\|\}\}$  for the two constants  $C_9$  and  $C_{10}$ , which means  $\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}$  is bounded from below. Thus, the first term tends to 0. Second, since the class  $\{I(V_s \leq t) \Delta_s / \mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \}_{v=V_s} : t \in [0, \tau]\}$  is also a Glivenko-Cantelli class, the second term vanishes as  $m_s$  goes to infinity.

Therefore, we conclude that  $\bar{\Lambda}_s(t)$  uniformly converges to

$$\mathbf{E} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v) Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right]_{v=V_s}. \tag{3.12}$$

We can easily verify that (3.12) is equal to  $\Lambda_{s0}(t)$ . Thus, the claim that  $\bar{\Lambda}_s(t)$  uniformly converges to  $\Lambda_{s0}(t)$  in  $[0, \tau]$  has been proved.

From the construction of  $\bar{\Lambda}_s(t)$ , we obtain that

$$\widehat{\Lambda}_s(t) = \int_0^t \frac{d\widehat{\Lambda}_s(v)}{d\bar{\Lambda}_s(v)} d\bar{\Lambda}_s(v) = \int_0^t \frac{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}} d\bar{\Lambda}_s(v). \quad (3.13)$$

$\widehat{\Lambda}_s(t)$  is absolutely continuous with respect to  $\bar{\Lambda}_s(t)$ . On the other hand, since both  $\{I(V_s \geq v) : v \in [0, \tau]\}$  and  $\mathcal{F}$  are Glivenko-Cantelli classes,  $\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau]\}$  is also a Glivenko-Cantelli class. Thus, we have

$$\begin{aligned} & \sup_{v \in [0, \tau]} |(\mathbf{P}_{m_s} - \mathbf{P})\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}| + \sup_{v \in [0, \tau]} |(\mathbf{P}_{m_s} - \mathbf{P})\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}| \\ & \longrightarrow 0 \quad \text{a.s.} \end{aligned}$$

By the bounded convergence theorem and the fact that  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}^*$  and  $\widehat{\Lambda}_s$  converges to  $\Lambda_s^*$ , for each  $v$ ,  $\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\} \longrightarrow \mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}$ ; moreover, it is straightforward to check the derivative of  $\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}$  with respect to  $v$ . Thus, by the Arzela-Ascoli theorem, uniformly in  $[0, \tau]$ ,

$$\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\} \longrightarrow \mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}.$$

Then, combining the above result and (3.13), it holds that, uniformly in  $[0, \tau]$ ,

$$\frac{\widehat{\Lambda}_s\{v\}}{\bar{\Lambda}_s\{v\}} = \frac{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}} \longrightarrow \frac{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}}. \quad (3.14)$$

After taking limits on both sides of (3.13), we obtain that

$$\Lambda_s^*(t) = \int_0^t \frac{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}} d\Lambda_{s0}(v), \quad (3.15)$$

Therefore, since  $\Lambda_{s0}(t)$  is differentiable with respect to the Lebesgue measure, so is

$\Lambda_s^*(t)$ ; that is, (3.15) is equal to

$$\int_0^t \frac{d\Lambda_s^*(v)}{d\Lambda_{s0}(v)} d\Lambda_{s0}(v). \quad (3.16)$$

And we denote  $\lambda_s^*(t)$  as the derivative of  $\Lambda_s^*(t)$ . Additionally, from (3.14) ~ (3.16), note that  $\widehat{\Lambda}_s\{V_s\}/\bar{\Lambda}_s\{V_s\}$  uniformly converges to  $d\Lambda_s^*(V_s)/d\Lambda_{s0}(V_s) = \lambda_s^*(V_s)/\lambda_{s0}(V_s)$ . Therefore, a second conclusion is that  $\widehat{\Lambda}_s$  uniformly converges to  $\Lambda_s^*$  since  $\Lambda_s^*$  is continuous.

On the other hand,

$$\begin{aligned} & n^{-1}l_n(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}) - n^{-1}l_n(\boldsymbol{\theta}_0, \bar{\boldsymbol{\Lambda}}) \\ &= \sum_{s=1}^S \left( \mathbf{P}_{m_s} \left[ \Delta_s \log \frac{\widehat{\Lambda}_s\{V_s\}}{\bar{\Lambda}_s\{V_s\}} \right] + \mathbf{P}_{m_s} \left[ \log \frac{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}}{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}} \right] \right) \\ &\geq 0. \end{aligned} \quad (3.17)$$

Using the result of Section 3.5.3.1 and similar arguments as above, we can verify that

$$\log \frac{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}}{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}}$$

belongs to a Glivenko-Cantelli class and

$$\mathbf{P} \left[ \log \frac{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}}{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}} \right] \rightarrow \mathbf{P} \left[ \log \frac{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}}{\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}} \right].$$

Since  $\widehat{\Lambda}_s\{V_s\}/\bar{\Lambda}_s\{V_s\}$  uniformly converges to  $\lambda_s^*\{V_s\}/\lambda_{s0}\{V_s\}$ , we obtain that, from (3.17),

$$\mathbf{P} \left[ \log \left\{ \frac{(\lambda_s^*(V_s))^{\Delta_s} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}}{(\lambda_{s0}(V_s))^{\Delta_s} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}} \right\} \right] \geq 0.$$

Note that the left-hand side of the inequality is the negative Kullback-Leibler informa-

tion. Then, the equality holds with probability one, and it immediately follows

$$(\lambda_s^*(V_s))^{\Delta_s} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} = (\lambda_{s0}(V_s))^{\Delta_s} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}. \quad (3.18)$$

Our proof will be completed if we can show  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$  from (3.18). Since (3.18) holds with probability one, (3.18) holds for any  $(V_s, \Delta_s = 1)$  and the case  $(V_s = \tau, \Delta_s = 0)$ , but may not hold for  $(V_s, \Delta_s = 0)$  when  $V \in (0, \tau)$ . However, we can show that (3.18) is also true for  $(V_s, \Delta_s = 0)$  when  $V_s \in (0, \tau)$ . To do this, treating both sides of (3.18) as functions of  $V_s$ , we integrate these functions over an interval  $(V_s, \tau)$  for  $\Delta_s = 0$  as the following;

$$\int_{V_s}^{\tau} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} = \int_{V_s}^{\tau} \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}$$

to obtain that

$$\begin{aligned} & \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0, V_s=\tau} - \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0, V_s=V_s} \\ &= \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0, V_s=\tau} - \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0, V_s=V_s}. \end{aligned}$$

After comparing this above equality with another following equality, which is given by (3.18) at  $\Delta_s = 0$  and  $V_s = \tau$ ,

$$\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0, V_s=\tau} = \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0, V_s=\tau},$$

we obtain

$$\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0, V_s=V_s} = \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0, V_s=V_s},$$

and therefore

$$\int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0} = \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0};$$

that is, (3.18) also holds for any  $V_s$  and  $\Delta_s = 0$ .

Thus, first to show that  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_0$ ,  $\boldsymbol{\phi}^* = \boldsymbol{\phi}_0$  and  $\boldsymbol{\Sigma}_b^* = \boldsymbol{\Sigma}_{b0}$ , we let  $\Delta_s = 0$  and  $V_s = 0$  in (3.18). After integrating over  $\mathbf{b}$ , we have that, with probability one,

$$\begin{aligned} & \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b} \Big|_{\Delta_s=0, V_s=0} = \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}, \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b} \Big|_{\Delta_s=0, V_s=0} \\ \Rightarrow & \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta}^* + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\boldsymbol{\beta}^*, \mathbf{b})}{A(D(t_j; \boldsymbol{\phi}^*))} + C(Y_j; D(t_j; \boldsymbol{\phi}^*)) \right] \right\} \\ & \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b^*|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_b^{*-1} \mathbf{b} \right\} d\mathbf{b} \\ = & \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\boldsymbol{\beta}_0, \mathbf{b})}{A(D(t_j; \boldsymbol{\phi}_0))} + C(Y_j; D(t_j; \boldsymbol{\phi}_0)) \right] \right\} \\ & \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_{b0}|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b} \right\} d\mathbf{b} \\ \Rightarrow & \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \boldsymbol{\phi}^*))} + C(Y_j; D(t_j; \boldsymbol{\phi}^*)) \right] \right\} \times |\boldsymbol{\Sigma}_b^*|^{-1/2} \\ & \quad \times \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j \mathbf{b}}{A(D(t_j; \boldsymbol{\phi}^*))} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b})}{A(D(t_j; \boldsymbol{\phi}^*))} - \frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_b^{*-1} \mathbf{b} \right\} d\mathbf{b} \\ = & \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}_0}{A(D(t_j; \boldsymbol{\phi}_0))} + C(Y_j; D(t_j; \boldsymbol{\phi}_0)) \right] \right\} \times |\boldsymbol{\Sigma}_{b0}|^{-1/2} \\ & \quad \times \int_{\mathbf{b}} \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j \mathbf{b}}{A(D(t_j; \boldsymbol{\phi}_0))} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \boldsymbol{\phi}_0))} - \frac{1}{2} \mathbf{b}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b} \right\} d\mathbf{b}. \end{aligned}$$

The left hand side becomes

$$\exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \boldsymbol{\phi}^*))} + C(Y_j; D(t_j; \boldsymbol{\phi}^*)) \right] \right\} \times |\boldsymbol{\Sigma}_b^*|^{-1/2}$$



$$\begin{aligned}
& \times \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \left[ (\Sigma_b^{*-1/2} \mathbf{b})^T (\Sigma_b^{*-1/2} \mathbf{b}) - 2 \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \mathbf{b} \right. \right. \\
& \quad \left. \left. + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right] \right. \\
& \quad \left. + \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b})}{A(D(t_j; \phi^*))} \right\} d\mathbf{b} \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right. \\
& \quad \left. + \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right\} \times |\Sigma_b^*|^{-1/2} \\
& \quad \times \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \left[ \Sigma_b^{*-1/2} \mathbf{b} - \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right]^T \left[ \Sigma_b^{*-1/2} \mathbf{b} - \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right] \right\} \\
& \quad \times \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b})}{A(D(t_j; \phi^*))} \right\} d\mathbf{b} \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right. \\
& \quad \left. + \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right\} \times \mathbb{E} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b})}{A(D(t_j; \phi^*))} \right\} \right].
\end{aligned} \tag{3.19}$$

Likewise, the right-hand side becomes

$$\begin{aligned}
& \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}_0}{A(D(t_j; \phi_0))} + C(Y_j; D(t_j; \phi_0)) \right] + \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T \right\} \\
& \quad \times \mathbb{E} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right].
\end{aligned} \tag{3.20}$$

Then, to compare the coefficients of  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{Y}$  in the exponential part and the constant term out of the exponential part from (3.19) and (3.20), we have

$$\left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T = \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T, \tag{3.21}$$

$$\sum_{j=1}^{n_N} \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} = \sum_{j=1}^{n_N} \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}_0}{A(D(t_j; \phi_0))}, \quad (3.22)$$

and

$$\mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b})}{A(D(t_j; \phi^*))} \right\} \right] = \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \quad (3.23)$$

Define  $\tilde{\mathbf{X}}_j^* = \tilde{\mathbf{X}}_j / A(D(t_j; \phi^*))$  and  $\tilde{\mathbf{X}}_{j0} = \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0))$  and  $\tilde{\mathbf{X}}^* = (\tilde{\mathbf{X}}_1^{*T}, \dots, \tilde{\mathbf{X}}_{n_N}^{*T})^T$  and  $\tilde{\mathbf{X}}_0 = (\tilde{\mathbf{X}}_{10}^{*T}, \dots, \tilde{\mathbf{X}}_{n_N 0}^{*T})^T$ . Then, (3.21) can be expressed as

$$\mathbf{Y}^T \tilde{\mathbf{X}}^* \tilde{\mathbf{X}}^{*T} \mathbf{Y} = \mathbf{Y}^T \tilde{\mathbf{X}}_0 \tilde{\mathbf{X}}_0^T \mathbf{Y},$$

and we obtain  $\tilde{\mathbf{X}}^* \tilde{\mathbf{X}}^{*T} = \tilde{\mathbf{X}}_0 \tilde{\mathbf{X}}_0^T$  for the coefficients of  $\mathbf{Y}^T \mathbf{Y}$ . For the  $j$ -th diagonal element, we have

$$\tilde{\mathbf{X}}_j^* \tilde{\mathbf{X}}_j^{*T} = \tilde{\mathbf{X}}_{j0} \tilde{\mathbf{X}}_{j0}^T \Rightarrow \frac{\tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \frac{\tilde{\mathbf{X}}_j^{*T}}{A(D(t_j; \phi^*))} = \frac{\tilde{\mathbf{X}}_{j0}}{A(D(t_j; \phi_0))} \frac{\tilde{\mathbf{X}}_{j0}^T}{A(D(t_j; \phi_0))}.$$

By assumption (A5),  $(A(D(t_j; \phi^*)))^2 = (A(D(t_j; \phi_0)))^2$ . Then, we obtain  $A(D(t_j; \phi^*)) = A(D(t_j; \phi_0))$  since both  $A(D(t_j; \phi^*))$  and  $A(D(t_j; \phi_0))$  are positive by the assumption for dispersion parameter of the generalized linear mixed model. By the continuous mapping theorem, we obtain  $D(t_j; \phi^*) = D(t_j; \phi_0)$ . By the similar argument, for the comparison of the coefficients of  $\mathbf{Y}$ , (3.22) can be written as

$$\mathbf{Y}^T \mathbf{X}^* \boldsymbol{\beta}^* = \mathbf{Y}^T \mathbf{X}_0 \boldsymbol{\beta}_0 \Rightarrow \mathbf{X}^* \boldsymbol{\beta}^* = \mathbf{X}_0 \boldsymbol{\beta}_0,$$

where the  $j$ -th elements  $(\mathbf{X}_j / A(D(t_j; \phi^*))) \boldsymbol{\beta}^* = (\mathbf{X}_j / A(D(t_j; \phi_0))) \boldsymbol{\beta}_0$ . By the result  $A(D(t_j; \phi^*)) = A(D(t_j; \phi_0))$  and assumption (A5), we obtain  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_0$ . In (3.23) for the constant term, note that the random effect  $\mathbf{b}$  on the left-hand side follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_b^{*1/2} \left( \sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi^*)) \right)^T$  and covariance  $\boldsymbol{\Sigma}_b^*$  while the random effect  $\mathbf{b}$  on the right-hand side follows a multivariate normal distribution

bution with mean  $\Sigma_{b0}^{1/2}(\sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0)))^T$  and covariance  $\Sigma_{b0}$ . Since  $\beta^* = \beta_0$  and  $\phi^* = \phi_0$ ,  $\sum_{j=1}^{n_N} B(\beta^*; \mathbf{b}) / A(D(t_j; \phi^*)) = \sum_{j=1}^{n_N} B(\beta_0; \mathbf{b}) / A(D(t_j; \phi_0))$ . Thus, to hold the equality of the expected values in (3.23), the random effects  $\mathbf{b}$  on both sides follow the same distribution; that is,  $\Sigma_b^* = \Sigma_{b0}$ .

Next, to show that  $\psi^* = \psi_0$ ,  $\gamma^* = \gamma_0$  and  $\Lambda_s^* = \Lambda_{s0}$ , we let  $\Delta_s = 0$  in (3.18). Through the similar arguments done for the proof of  $\beta^* = \beta_0$ ,  $\phi^* = \phi_0$  and  $\Sigma_b^* = \Sigma_{b0}$ , we obtain

$$\begin{aligned} & \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta^*; \mathbf{b})}{A(D(t_j; \phi^*))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}) + \mathbf{Z}(t)\gamma^* \} d\Lambda_s^*(t) \right\} \right] \\ &= \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t)\gamma_0 \} d\Lambda_{s0}(t) \right\} \right], \quad (3.24) \end{aligned}$$

where the random effects  $\mathbf{b}$  follow a multivariate normal distribution with mean  $\Sigma_{b0}^{1/2}(\sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0)))^T$  and covariance  $\Sigma_{b0}$ . For any fixed  $\tilde{\mathbf{X}}$ , treating  $\tilde{\mathbf{X}}^T \mathbf{Y}$  as a parameter in this normal family,  $\mathbf{b}$  is the complete statistic for  $\tilde{\mathbf{X}}^T \mathbf{Y}$ . Therefore,

$$\begin{aligned} & \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta^*; \mathbf{b})}{A(D(t_j; \phi^*))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}) + \mathbf{Z}(t)\gamma^* \} d\Lambda_s^*(t) \right\} \\ &= \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t)\gamma_0 \} d\Lambda_{s0}(t) \right\}. \end{aligned}$$

Since  $\beta^* = \beta_0$  and  $\phi^* = \phi_0$ , equivalently, we have

$$\exp \{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}) + \mathbf{Z}(t)\gamma^* \} \lambda_s^*(t) = \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t)\gamma_0 \} \lambda_{s0}(t).$$

By assumptions (A2) and (A5),  $\psi^* = \psi_0$ ,  $\gamma^* = \gamma_0$  and  $\Lambda_s^* = \Lambda_{s0}$ .

Since all the three steps are completed, we can conclude that, with probability one,  $\hat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\hat{\boldsymbol{\Lambda}}$  converges to  $\boldsymbol{\Lambda}_0$  in  $[0, \tau]$ . Moreover, as mentioned in the beginning of this proof for consistency, since  $\boldsymbol{\Lambda}_0$  is continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\hat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$  almost

surely. Therefore, Theorem 3.1 is proved.

### 3.5.2 Proof of asymptotic normality

Asymptotic distribution for the proposed estimator can be shown if we can verify the conditions of Theorem 3.3.1 (p310) in van der Vaart and Wellner (1996). Then, we will show that the distribution is normal. For completeness, we state this theorem below following Theorem 4 in Appendix A of Parner (1998).

**Theorem 3.3.** (Theorem 3.3.1 in van der Vaart and Wellner, 1996) *Let  $U_n$  and  $U$  be random maps and a fixed map, respectively, from  $\xi$  to a Banach space such that:*

- (a)  $\sqrt{n}(U_n - U)(\widehat{\xi}_n) - \sqrt{n}(U_n - U)(\xi_0) = o_P^*(1 + \sqrt{n}\|\widehat{\xi}_n - \xi_0\|)$ .
- (b) *The sequence  $\sqrt{n}(U_n - U)(\xi_0)$  converges in distribution to a tight random element  $\mathbf{W}$ .*
- (c) *the function  $\xi \rightarrow U(\xi)$  is Fréchet differentiable at  $\xi_0$  with a continuously invertible derivative  $\nabla U_{\xi_0}$  (on its range).*
- (d)  $U_{\xi_0}$  and  $\widehat{\xi}_n$  satisfies  $U_n(\widehat{\xi}_n) = o_P^*(n^{-1/2})$  and converges in outer probability to  $\xi_0$ .

Then  $\sqrt{n}(\widehat{\xi}_n - \xi_0) \Rightarrow \nabla U_{\xi_0}^{-1} \mathbf{W}$ .

We will prove the conditions (a)~(d). In our situation, the parameter  $\xi_s = (\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) \in \Xi = \{(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, s = 1, \dots, S\}$  for a fixed small constant  $\delta$ . We note that  $\Xi$  is a convex set. Define a set  $\mathcal{H} = \{(\mathbf{h}_1, h_2) : \|\mathbf{h}_1\| \leq 1, \|h_2\|_V \leq 1\}$ , where  $\|h_2\|_V$  is the total variation of  $h_2$  in  $[0, \tau]$  defined as

$$\sup_{0=t_0 \leq t_2 \leq \dots \leq t_k = \tau} \sum_{j=1}^k |h_2(t_j) - h_2(t_{j-1})|.$$

Furthermore, we define that, for stratum  $s$ ,

$$U_{m_s}(\xi_s)(\mathbf{h}_1, h_2) = \mathbf{P}_{m_s}\{l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2]\}$$

and

$$U_s(\xi_s)(\mathbf{h}_1, h_2) = \mathbf{P}\{l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2]\},$$

where  $l_\theta(\boldsymbol{\theta}, \Lambda_s)$  is the first derivative of the log-likelihood function from one single subject belonging to stratum  $s$ , denoted by  $l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$ , with respect to  $\boldsymbol{\theta}$ , and  $l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)$  is the derivative of  $l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_{s\varepsilon})$  at  $\varepsilon = 0$ , where  $\Lambda_{s\varepsilon}(t) = \int_0^t (1 + \varepsilon h_2(u)) d\Lambda_{s0}(u)$ . Therefore, we can see that both  $U_{m_s}$  and  $U_s$  map from  $\Xi$  to  $\ell^\infty(\mathcal{H})$  and  $\sqrt{m_s}\{U_{m_s}(\xi_s) - U_s(\xi_s)\}$  is an empirical process in the space  $\ell^\infty(\mathcal{H})$ .

Denote  $(\mathbf{h}_1^\beta, \mathbf{h}_1^\phi, \mathbf{h}_1^b, \mathbf{h}_1^\psi, \mathbf{h}_1^\gamma)$  as the corresponding components of  $\mathbf{h}_1$  for the parameters  $(\boldsymbol{\beta}, \boldsymbol{\phi}, \text{Vec}(\boldsymbol{\Sigma}_b), \boldsymbol{\psi}, \boldsymbol{\gamma})$ , respectively. From Section 3.5.3.2, for any  $(\mathbf{h}_1, h_2) \in \mathcal{H}$ , the class

$$\begin{aligned} \mathcal{G} = & \left\{ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2], \right. \\ & \left. \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, (\mathbf{h}_1, h_2) \in \mathcal{H} \right\} \end{aligned}$$

is shown as P-Donsker (Section 2.1 of van der Vaart and Wellner (1996)), and it is also implied that

$$\sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \mathbf{P} \left[ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right]^2 \longrightarrow 0$$

as  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \rightarrow 0$ . Then, we conclude the followings:

- (a) follows from Lemma 3.3.5 (p311) of van der Vaart and Wellner (1996).
- (b) holds as a result of Section 3.5.3.2 and the convergence is defined in the metric space  $\ell^\infty(\mathcal{H})$  by the Donsker theorem (Section 2.5 of van der Vaart and Wellner

(1996).

(d) is true because  $(\widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)$  maximizes  $\mathbf{P}_{m_s} l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$ ,  $(\boldsymbol{\theta}_0, \Lambda_{s0})$  maximizes  $\mathbf{P} l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$ , and  $(\widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)$  converges to  $(\boldsymbol{\theta}_0, \Lambda_{s0})$  from Theorem 3.1.

Now, we need to verify the conditions in (c). Since the proof of the first half in (c), that the function  $\xi \rightarrow U(\xi)$  is *Fréchet* differentiable at  $\xi_0$ , is given in Section 3.5.3.3, we will only prove that the derivative  $\nabla U_{\xi_0}$  is continuously invertible on its range  $\ell^\infty(\mathcal{H})$ . According to Section 3.5.3.3,  $\nabla U_{\xi_0}$  can be expressed as follows: for any  $(\boldsymbol{\theta}_1, \Lambda_{s1})$  and  $(\boldsymbol{\theta}_2, \Lambda_{s2})$  in  $\Xi$ ,

$$\nabla U_{\xi_0}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\Lambda_{s1} - \Lambda_{s2})(t), \quad (3.25)$$

where both  $\Omega_1$  and  $\Omega_2$  are linear operators on  $\mathcal{H}$ , and  $\Omega = (\Omega_1, \Omega_2)$  maps  $\mathcal{H} \subset \mathbf{R}^d \times \text{BV}[0, \tau]$  to  $\mathbf{R}^d \times \text{BV}[0, \tau]$ , where  $\text{BV}[0, \tau]$  contains all the functions with finite total variation in  $[0, \tau]$ . The explicit expressions of  $\Omega_1$  and  $\Omega_2$  are given in Section 3.5.3.3. From (3.25), we can treat  $(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})$  as an element in  $\ell^\infty(\mathcal{H})$  via the following definition:

$$(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\Lambda_{s1} - \Lambda_{s2})(t), \quad \forall (\mathbf{h}_1, h_2) \in \mathbf{R}^d \times \text{BV}[0, \tau].$$

Then  $\nabla U_{\xi_0}$  can be expanded as a linear operator from  $\ell^\infty(\mathcal{H})$  to itself. Therefore, if we can show that there exists some positive constant  $\varepsilon$  such that  $\varepsilon \mathcal{H} \subset \Omega(\mathcal{H})$ , then we will have that for any  $(\delta \boldsymbol{\theta}, \delta \Lambda_s) \in \ell^\infty(\mathcal{H})$ ,

$$\begin{aligned} \|\nabla U_{\xi_0}(\delta \boldsymbol{\theta}, \delta \Lambda_s)\|_{\ell^\infty(\mathcal{H})} &= \sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \left| \delta \boldsymbol{\theta}^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d\delta \Lambda_s(t) \right| \\ &= \|(\delta \boldsymbol{\theta}, \delta \Lambda_s)\|_{\ell^\infty(\Omega(\mathcal{H}))} \geq \varepsilon \|(\delta \boldsymbol{\theta}, \delta \Lambda_s)\|_{\ell^\infty(\mathcal{H})}, \end{aligned}$$

and  $\nabla U_{\xi_0}$  will be continuously invertible.

Note that to prove  $\varepsilon \mathcal{H} \subset \Omega(\mathcal{H})$  for some  $\varepsilon$  is equivalent to showing that  $\Omega$  is invertible. We also note from Section 3.5.3.3, that  $\Omega$  is the summation of an invertible operator and a compact operator. By Theorem 4.25 of Rudin (1973), for the proof of the invertibility of  $\Omega$ , it is sufficient to verify that  $\Omega$  is one to one: if  $\Omega[\mathbf{h}_1, h_2] = 0$ , then, by choosing  $\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 = \varepsilon^* \mathbf{h}_1$  and  $\Lambda_{s1} - \Lambda_{s2} = \varepsilon^* \int h_2 d\Lambda_{s0}$  in (3.25) for a small constant  $\varepsilon^*$ , we obtain

$$\nabla U_{\xi_0}(\mathbf{h}_1, \int h_2 d\Lambda_{s0})[\mathbf{h}_1, h_2] = \varepsilon^* (\mathbf{h}_1^T, h_2) \begin{pmatrix} \Omega_1[\mathbf{h}_1, h_2] \\ \Omega_2[\mathbf{h}_1, h_2] \end{pmatrix} = \varepsilon^* (\mathbf{h}_1^T, h_2) \Omega[\mathbf{h}_1, h_2] = 0.$$

By the definition of  $\nabla U_{\xi_0}$ , we note that  $\nabla U_{\xi_0}(\mathbf{h}_1, \int h_2 d\Lambda_{s0})[\mathbf{h}_1, h_2]$  is the negative information matrix in the submodel  $(\boldsymbol{\theta}_0 + \varepsilon \mathbf{h}_1, \Lambda_{s0} + \varepsilon \int h_2 d\Lambda_{s0})$ . Thus, the score function along this submodel should be zero with probability one; that is,  $l_{\boldsymbol{\theta}}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] = 0$ ; that is, with probability one, for the numerator of the score function

$$\begin{aligned} 0 = & \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \frac{\mathbf{b}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right. \\ & + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))^2} \right) (A'(D(t_j; \phi_0))) h_1^\phi + C'(Y_j; D(t_j; \phi_0)) h_1^\phi \right\} \\ & + \sum_{j=1}^{n_N} \left\{ \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi_0))} \mathbf{h}_1^\beta - B'(\boldsymbol{\beta}_0; \mathbf{b}) \mathbf{h}_1^\beta \right\} + \Delta_s \{ (\tilde{\mathbf{Z}}(V_s) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(V_s) \mathbf{h}_1^\gamma \} \\ & - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t) \mathbf{h}_1^\gamma \} d\Lambda_{s0}(t) \Big] d\mathbf{b} \\ & + \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \Delta_s h_2(V_s) - \int_0^{V_s} h_2(t) \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \} d\Lambda_{s0}(t) \right] d\mathbf{b}, \end{aligned} \tag{3.26}$$

where  $A'(D(t_j; \phi_0))$  and  $C'(Y_j; D(t_j; \phi_0))$  are the derivatives of  $A(\phi_j)$  and  $C(Y_j; \phi_j)$  with

respect to  $\phi_j$  evaluated at  $D(t_j; \phi_0)$  and  $B'(\beta_0; \mathbf{b})$  is the derivative of  $B(\beta; \mathbf{b})$  with respect to  $\beta$  evaluated at  $\beta_0$ . Note that (3.26) holds with probability one, so it may not hold for any  $V_s \in [0, \tau]$  when  $\Delta_s = 0$ . However, by the similar arguments done in Section 3.5.1, if we integrate both sides from  $V_s$  to  $\tau$  and subtract the obtained equation from (3.26) at  $\Delta_s = 0$  and  $V_s = \tau$ , it is easily shown that (3.26) also holds for any  $V_s \in [0, \tau]$  when  $\Delta_s = 0$ . Hence, the proof of the invertibility of  $\Omega$  will be completed if we can show  $\mathbf{h}_1 = 0$  and  $h_2(t) = 0$  from (3.26).

To show  $\mathbf{h}_1 = 0$ , particularly we let  $\Delta_s = 0$  and  $V_s = 0$  in (3.26) and obtain

$$\begin{aligned}
0 &= \int \mathbf{b} G(\mathbf{b}, \mathbf{O}; \theta_0, \Lambda_{s0}) \times \left[ \frac{\mathbf{b}^T \Sigma_{b0}^{-1} \mathbf{D}_b \Sigma_{b0}^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\Sigma_{b0}^{-1} \mathbf{D}_b) \right. \\
&\quad + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j(\mathbf{X}_j \beta_0 + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))^2} \right) (A'(D(t_j; \phi_0))) \mathbf{h}_1^\phi + C'(Y_j; D(t_j; \phi_0)) \mathbf{h}_1^\phi \right\} \\
&\quad + \sum_{j=1}^{n_N} \left\{ \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi_0))} \mathbf{h}_1^\beta - B'(\beta_0; \mathbf{b}) \mathbf{h}_1^\beta \right\} \Big] d\mathbf{b} \\
&= \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \frac{\mathbf{b}^T \Sigma_{b0}^{-1} \mathbf{D}_b \Sigma_{b0}^{-1} \mathbf{b}}{2} \right] \\
&\quad + \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \left\{ - \frac{1}{2} \text{Tr}(\Sigma_{b0}^{-1} \mathbf{D}_b) \right. \right. \\
&\quad \quad + \sum_{j=1}^{n_N} \left( - \left( \frac{Y_j \mathbf{X}_j \beta_0}{(A(D(t_j; \phi_0)))^2} \right) (A'(D(t_j; \phi_0))) \mathbf{h}_1^\phi + C'(Y_j; D(t_j; \phi_0)) \mathbf{h}_1^\phi \right) \\
&\quad \quad + \sum_{j=1}^{n_N} \left( - \left( \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi_0))} \right) \mathbf{h}_1^\beta \right) \Big\} \\
&\quad \left. + \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j \tilde{\mathbf{X}}_j \mathbf{b} - B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))^2} \right) (A'(D(t_j; \phi_0))) \mathbf{h}_1^\phi \right. \right. \right. \\
&\quad \quad \quad \left. \left. \left. - B'(\beta_0; \mathbf{b}) \mathbf{h}_1^\beta \right\} \right] \right]. \tag{3.27}
\end{aligned}$$

We first examine the coefficient for  $\mathbf{Y}$  in (3.27).



$$\begin{aligned}
& \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \\
& \quad \times \sum_{j=1}^{n_N} \left\{ - \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi_0))} \left( \frac{\boldsymbol{\beta}_0}{A(D(t_j; \phi_0))} A'(D(t_j; \phi_0)) \mathbf{h}_1^\phi - \mathbf{h}_1^\beta \right) \right\} \\
& \quad + \sum_{j=1}^{n_N} \left\{ - \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))^2} (A'(D(t_j; \phi_0))) \right\} \mathbf{h}_1^\phi \mathbb{E} \left[ \mathbf{b} \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \\
& = \sum_{j=1}^{n_N} \left\{ - \frac{Y_j}{A(D(t_j; \phi_0))} \left[ \mathbf{X}_j \left( \frac{\boldsymbol{\beta}_0}{A(D(t_j; \phi_0))} (A'(D(t_j; \phi_0))) \mathbf{h}_1^\phi - \mathbf{h}_1^\beta \right) \right. \right. \\
& \quad \times \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \\
& \quad \left. \left. + \frac{\tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} (A'(D(t_j; \phi_0))) \mathbf{h}_1^\phi \mathbb{E} \left[ \mathbf{b} \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \right] \right\} \\
& = \sum_{j=1}^{n_N} \left\{ - \frac{Y_j}{A(D(t_j; \phi_0))} \left[ \left[ \mathbf{X}_j \left( \frac{\boldsymbol{\beta}_0}{A(D(t_j; \phi_0))} (A'(D(t_j; \phi_0))) \right) \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \right. \right. \right. \\
& \quad \left. \left. + \frac{\tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} (A'(D(t_j; \phi_0))) \mathbb{E} \left[ \mathbf{b} \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \right] \mathbf{h}_1^\phi \right. \right. \\
& \quad \left. \left. - \left[ \mathbf{X}_j \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \right] \right] \mathbf{h}_1^\beta \right] \right\} \\
& = 0
\end{aligned}$$

Based on assumption (A5),  $\mathbf{h}_1^\phi = 0$  and  $\mathbf{h}_1^\beta = 0$ .

Then, we examine the constant terms without  $\mathbf{Y}$  in (3.27). Since  $\mathbf{h}_1^\phi = 0$  and  $\mathbf{h}_1^\beta = 0$ , (3.27) becomes

$$\begin{aligned}
& \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \frac{\mathbf{b}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}}{2} \right] \\
& \quad + \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \left\{ - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right\} \right] \\
& = \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \left\{ \frac{\mathbf{b}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right\} \right] = 0,
\end{aligned}$$

where  $\mathbf{b}$  follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_{b0}^{1/2} \left[ \sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{Z}}_j / A(D(t_j; \phi_0))) \right]$

$\phi_0)))]^T$  and covariance  $\Sigma_{b0}$ . For any fixed  $\tilde{\mathbf{X}}$ , treating  $\mathbf{X}^T \mathbf{Y}$  as a parameter in this normal family,  $\mathbf{b}$  is the complete statistic for  $\mathbf{X}^T \mathbf{Y}$ , therefore,

$$\exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} \right\} \times \left\{ \frac{\mathbf{b}^T \Sigma_{b0}^{-1} \mathbf{D}_b \Sigma_{b0}^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\Sigma_{b0}^{-1} \mathbf{D}_b) \right\} = 0.$$

Since  $\exp \left\{ - \sum_{j=1}^{n_N} (B(\beta_0; \mathbf{b})/A(D(t_j; \phi_0))) \right\} \neq 0$ ,  $[\mathbf{b}^T \Sigma_{b0}^{-1} \mathbf{D}_b \Sigma_{b0}^{-1} \mathbf{b} - \text{Tr}(\Sigma_{b0}^{-1} \mathbf{D}_b)]/2 = 0$  by (A5). Then, since  $\Sigma_{b0}^{-1} \neq 0$ , then  $\mathbf{D}_b = 0$ .

Next, we let  $\Delta_s = 0$  in (3.26) and obtain

$$\begin{aligned} 0 &= \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \theta_0, \Lambda_{s0}) \times \left[ \frac{\mathbf{b}^T \Sigma_{b0}^{-1} \mathbf{D}_b \Sigma_{b0}^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\Sigma_{b0}^{-1} \mathbf{D}_b) \right. \\ &\quad \left. + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j(\mathbf{X}_j \beta_0 + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))^2} \right) (A'(D(t_j; \phi_0))) h_1^\phi + C'(Y_j; D(t_j; \phi_0)) h_1^\phi \right\} \right. \\ &\quad \left. + \sum_{j=1}^{n_N} \left\{ \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi_0))} \mathbf{h}_1^\beta - B'(\beta_0; \mathbf{b}) \mathbf{h}_1^\beta \right\} \right. \\ &\quad \left. - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t) \gamma_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t) \mathbf{h}_1^\gamma \} d\Lambda_{s0}(t) \right] d\mathbf{b} \\ &\quad + \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \theta_0, \Lambda_{s0}) \times \left[ - \int_0^{V_s} h_2(t) \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t) \gamma_0 \} d\Lambda_{s0}(t) \right] d\mathbf{b} \\ &= \mathbb{E} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b})}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t) \gamma_0 \} d\Lambda_{s0}(t) \right\} \right. \\ &\quad \left. \times \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}) + \mathbf{Z}(t) \gamma_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t) \mathbf{h}_1^\gamma + h_2(t) \} d\Lambda_{s0}(t) \right], \end{aligned} \tag{3.28}$$

where  $\mathbf{b}$  follows a multivariate normal distribution with mean  $\Sigma_{b0}^{1/2} [\sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{Z}}_j / A(D(t_j; \phi_0)))]^T$  and covariance  $\Sigma_{b0}$ . Likewise, for any fixed  $\tilde{\mathbf{X}}$ , treating  $\mathbf{X}^T \mathbf{Y}$  as a parameter in this normal family,  $\mathbf{b}$  is the complete statistic for  $\mathbf{X}^T \mathbf{Y}$  and therefore

$$\begin{aligned} & \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b})}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} d\Lambda_{s0}(t) \right\} \\ & \times \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t)\mathbf{h}_1^\gamma + h_2(t) \} d\Lambda_{s0}(t) = 0. \end{aligned}$$

Since  $\exp \{ - \sum_{j=1}^{n_N} (B(\boldsymbol{\beta}_0; \mathbf{b})/A(D(t_j; \phi_0))) - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} d\Lambda_{s0}(t) \} \neq 0$ , equivalently

$$\int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t)\mathbf{h}_1^\gamma + h_2(t) \} d\Lambda_{s0}(t) = 0$$

by assumption (A5). From assumption (A5), this immediately gives  $\mathbf{h}_1^\psi = 0$ ,  $\mathbf{h}_1^\gamma = 0$  and  $h_2(t) = 0$ . Hence, the proof of condition (c) is completed.

Since the conditions (a)–(d) have been proved, Theorem 3.3.1 of van der Vaart and Wellner (1996) concludes that  $\sqrt{m_s}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})$  weakly converges to a tight random element in  $\ell^\infty(\mathcal{H})$ . Furthermore, we obtain

$$\begin{aligned} & \sqrt{m_s} \nabla U_{\psi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\ & = \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P}) \{ l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \} + o_P(1), \end{aligned} \quad (3.29)$$

where  $o_P(1)$  is a random variable which converges to zero in probability in  $\ell^\infty(\mathcal{H})$ .

On the other hand, from (3.25), we have

$$\begin{aligned} & \sqrt{m_s} \nabla U_{\psi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\ & = \sqrt{m_s} \{ (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\widehat{\Lambda}_s - \Lambda_{s0})(t) \}. \end{aligned} \quad (3.30)$$

By denoting  $(\mathbf{h}_1^*, h_2^*) = \Omega^{-1}(\mathbf{h}_1, h_2)$ , we have  $(\mathbf{h}_1, h_2) = \Omega(\mathbf{h}_1^*, h_2^*)$ , and replacing  $(\mathbf{h}_1, h_2)$

with  $(\mathbf{h}_1^*, h_2^*)$  in (3.29) and (3.30) leads to the followings, respectively.

$$\begin{aligned} & \sqrt{m_s} \nabla U_{\psi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1^*, h_2^*] \\ &= \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P})\{l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]\} + o_P(1), \end{aligned}$$

and

$$\begin{aligned} & \sqrt{m_s} \nabla U_{\psi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1^*, h_2^*] \\ &= \sqrt{m_s} \left\{ (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1^*, h_2^*] + \int_0^\tau \Omega_2[\mathbf{h}_1^*, h_2^*] d(\widehat{\Lambda}_s - \Lambda_{s0})(t) \right\} \\ &= \sqrt{m_s} \left\{ (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\widehat{\Lambda}_s - \Lambda_{s0})(t) \right\}. \end{aligned}$$

Thus, we obtain

$$\begin{aligned} & \sqrt{m_s} \left\{ (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\widehat{\Lambda}_s - \Lambda_{s0})(t) \right\} \\ &= \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P})\{l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]\} + o_P(1). \quad (3.31) \end{aligned}$$

Note that the first term on the right-hand side in (3.31) is  $\sqrt{m_s}\{U_{m_s}(\boldsymbol{\theta}_0, \Lambda_{s0}) - U_s(\boldsymbol{\theta}_0, \Lambda_{s0})\}$ , which is an empirical process in the space  $\ell^\infty(\mathcal{H})$ , and it is shown that  $\mathcal{G}$  is P-Donsker in Section 3.5.3.2. Therefore,  $\sqrt{m_s}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})$  weakly converges to a Gaussian process in  $\ell^\infty(\mathcal{H})$ .

In particular, if we choose  $h_2 = 0$  in (3.31), then  $\widehat{\boldsymbol{\theta}}^T \mathbf{h}_1$  is a asymptotic linear estimator for  $\boldsymbol{\theta}_0^T \mathbf{h}_1$  with influence function being  $l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]$ . Since this influence function is in the linear space spanned by the score functions for  $\boldsymbol{\theta}_0$  and  $\Lambda_{s0}$ , Proposition 3.3.1 (p65) in Bickel, Klaassen, Ritov and Wellner (1993) concludes that the influence function is the same as the efficient influence function for  $\boldsymbol{\theta}_0^T \mathbf{h}_1$ ; that is  $\widehat{\boldsymbol{\theta}}$  is an efficient estimator for  $\boldsymbol{\theta}_0$  and Theorem 3.2 is proved.

### 3.5.3 Supplementary proofs

The proofs for P-Donsker property of the classes  $\mathcal{F}$  and  $\mathcal{G}$  needed in Sections 3.5.1 and 3.5.2 are presented in Sections 3.5.3.1~3.5.3.2 respectively. In Section 3.5.3.3, we prove *Fréchet* differentiability of  $U(\xi)$  at  $\xi_0$  and derive the derivative operator  $\nabla U_{\xi_0}$  use in Section 3.5.2.

#### 3.5.3.1 Proof of P-Donsker property of $\mathcal{F}$

We defined that a class  $\mathcal{F} = \{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau], \boldsymbol{\theta} \in \Theta, \Lambda_s \in \mathcal{A}, s = 1, \dots, S\}$ , where  $\mathcal{A} = \{\Lambda_s \in \mathbb{W}, \Lambda_s(0) = 0, \Lambda_s(\tau) \leq B_{s0}, s = 1, \dots, S\}$ ,  $B_{s0}$  is the constant given in the second step and  $\mathbb{W}$  contains all nondecreasing functions in  $[0, \tau]$ . We can rewrite  $Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$  as

$$Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) = Q_1(v, \mathbf{O}; \boldsymbol{\theta}) \frac{Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)}{Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)},$$

where

$$\begin{aligned} Q_1(v, \mathbf{O}; \boldsymbol{\theta}) &= \exp \left\{ \mathbf{Z}(v) \boldsymbol{\gamma} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T \right. \\ &\quad \left. + \frac{1}{2} (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T \right\}, \\ Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) &= \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\ &\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d\Lambda_s(t) \right\} d\mathbf{b}, \\ Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) &= \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_2(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\ &\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T \right\} d\Lambda_s(t) \right\} d\mathbf{b}, \end{aligned}$$

$$\begin{aligned} R(t) &= (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T, \quad B_1(\boldsymbol{\beta}; \mathbf{b}) = B(\boldsymbol{\beta}; g_1(\mathbf{b})), \quad B_2(\boldsymbol{\beta}; \mathbf{b}) = B(\boldsymbol{\beta}; g_2(\mathbf{b})), \\ g_1(\mathbf{b}) &= \boldsymbol{\Sigma}_b^{1/2} [\mathbf{b} + (\sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi))) + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T))^T] \text{ and } g_2(\mathbf{b}) = \boldsymbol{\Sigma}_b^{1/2} [\mathbf{b} + \end{aligned}$$

$$\left( \sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi))) + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right)^T \Big].$$

Using assumption (A2), we can easily show that  $Q_1(v, \mathbf{O}; \boldsymbol{\theta})$  is continuously differentiable with respect to  $v$  and  $\boldsymbol{\theta}$ , and

$$\|\nabla_{\boldsymbol{\theta}} Q_1(v, \mathbf{O}; \boldsymbol{\theta})\| + \left| \frac{d}{dv} Q_1(v, \mathbf{O}; \boldsymbol{\theta}) \right| \leq e^{k_1 + k_2 \|\mathbf{Y}\|}$$

for some positive constants  $k_1$  and  $k_2$ . Furthermore, it holds that

$$\begin{aligned} & \|\nabla_{\boldsymbol{\theta}} Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \\ & \leq \int_{\mathbf{b}} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times e^{k_3 \|\mathbf{b}\| + k_4 \|\mathbf{Y}\| + k_5} \times B_{s0} \right] d\mathbf{b} \\ & \leq e^{k_6 + k_7 \|\mathbf{Y}\|} \end{aligned}$$

$$\text{and} \quad \|\nabla_{\boldsymbol{\theta}} Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \leq e^{k_8 + k_9 \|\mathbf{Y}\|}$$

for some positive constants  $k_3, k_4, \dots, k_9$ . Additionally, note that, for any  $0 < \Lambda < \infty$ ,  $0 < e^{-\Lambda} < 1$  and  $e^{-\Lambda} < \Lambda$  and thus  $e^{-\Lambda_1} - e^{-\Lambda_2} < \Lambda_1 - \Lambda_2$  for any  $\Lambda_1$  and  $\Lambda_2$  over  $(0, \infty)$ . Hence,

$$\begin{aligned} & |Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s1}) - Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s2})| \\ & = \left| \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times \left[ \exp \left\{ -\int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} \right. \right. \right. \right. \\ & \quad \left. \left. \left. + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d\Lambda_{s1}(t) \right\} \right. \\ & \quad \left. \left. - \exp \left\{ -\int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{W}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} \right. \right. \right. \right. \right. \\ & \quad \left. \left. \left. \left. + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d\Lambda_{s2}(t) \right\} \right] d\mathbf{b} \right| \\ & \leq \left| \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \end{aligned}$$

$$\begin{aligned}
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \Big\} d(\Lambda_{s1} - \Lambda_{s2})(t) \Big] d\mathbf{b} \Big| \\
& = \left| \int_{\mathbf{b}} \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times \left[ \int_0^{V_s} \exp \left\{ - \frac{1}{2} \left[ \mathbf{b}^T \mathbf{b} - 2(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} \right. \right. \right. \right. \\
& \quad \left. \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T - (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right] \right. \right. \\
& \quad \left. \left. + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d(\Lambda_{s1} - \Lambda_{s2})(t) \right] d\mathbf{b} \Big| \\
& = \left| \int_{\mathbf{b}} \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times (2\pi)^{d_b/2} \right. \\
& \quad \times (2\pi)^{-d_b/2} \times \left[ \int_0^{V_s} \exp \left\{ - \frac{1}{2} \left[ \mathbf{b} - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right]^T \left[ \mathbf{b} - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right] \right\} \right. \\
& \quad \times \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \\
& \quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d(\Lambda_{s1} - \Lambda_{s2})(t) \right] d\mathbf{b} \Big| \\
& \leq \left| \mathbb{E}_b \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \right] (2\pi)^{d_b/2} \int_0^{V_s} \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \right. \\
& \quad \left. \left. + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d(\Lambda_{s1} - \Lambda_{s2})(t) \right] \Big| \\
& = K_0 \left| \int_0^{V_s} \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \right. \\
& \quad \left. \left. + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d(\Lambda_{s1} - \Lambda_{s2})(t) \right| \\
& = K_0 \left| - \int_0^{V_s} (\Lambda_{s1}(t) - \Lambda_{s2}(t)) \frac{d}{dt} \left[ \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \right. \right. \\
& \quad \left. \left. + \mathbf{Z}(t) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} \right] d(t) \right. \\
& \quad \left. + (\Lambda_{s1}(V_s) - \Lambda_{s2}(V_s)) \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \right. \\
& \quad \left. \left. + \mathbf{Z}(V_s) \boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \right)^T + R(V_s) \right\} \right| \\
& \leq K_0 \left[ \int_0^{V_s} |\Lambda_{s1}(t) - \Lambda_{s2}(t)| \left| \frac{d}{dt} \left[ \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \right. \right. \right. \right.
\end{aligned}$$

$$\begin{aligned}
& + \mathbf{Z}(t)\boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \Big\} \Big\| dt \\
& + \left| \Lambda_{s1}(V_s) - \Lambda_{s2}(V_s) \right| \exp \left\{ \frac{1}{2} (\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} ((\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T \right. \\
& \quad \left. + \mathbf{Z}(V_s)\boldsymbol{\gamma} + (\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(V_s) \circ \boldsymbol{\psi}^T) \right)^T + R(V_s) \right\} \\
& \leq e^{k_{10}+k_{11}\|\mathbf{Y}\|} \left\{ \left| \Lambda_{s1}(V_s) - \Lambda_{s2}(V_s) \right| + \int_0^\tau \left| \Lambda_{s1}(t) - \Lambda_{s2}(t) \right| dt \right\},
\end{aligned}$$

where  $K_0 = \mathbb{E}_b \left[ \exp \left\{ - \sum_{j=1}^{n_N} (B_1(\boldsymbol{\beta}; \mathbf{b}) / A(D(t_j; \phi))) \right\} \right] (2\pi)^{d_b/2}$ ,  $k_{10}$  and  $k_{11}$  are positive constants. Similarly,

$$\begin{aligned}
& |Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s1}) - Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s2})| \\
& \leq e^{k_{12}+k_{13}\|\mathbf{Y}\|} \left\{ \left| \Lambda_{s1}(V_s) - \Lambda_{s2}(V_s) \right| + \int_0^\tau \left| \Lambda_{s1}(t) - \Lambda_{s2}(t) \right| dt \right\},
\end{aligned}$$

where  $k_{12}$  and  $k_{13}$  are positive constants.

On the other hand, there exist positive constants  $k_{14}, \dots, k_{24}$  such that

$$\begin{aligned}
& |Q_1(v, \mathbf{O}; \boldsymbol{\theta})| \\
& = \left| \exp \left\{ \mathbf{Z}(v)\boldsymbol{\gamma} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T \right. \right. \\
& \quad \left. \left. + \frac{1}{2} (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T \right\} \right| \\
& \leq e^{k_{14}+k_{15}\|\mathbf{Y}\|}, \\
& |Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)| \\
& = \left| \int_{\mathbf{b}} \exp \left\{ - \frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t)\boldsymbol{\gamma} \right. \right. \right. \\
& \quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \right\} d\Lambda_s(t) \right\} d\mathbf{b} \right| \\
& \leq \left| \int_{\mathbf{b}} \exp \left\{ - \frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times \left[ 2 \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t)\boldsymbol{\gamma} \right. \right. \right.
\end{aligned}$$



$$\begin{aligned}
& +(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + R(t) \Big] d\Lambda_s(t) \Big] d\mathbf{b} \Big| \\
& \leq \left| \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} \right\} \times 2 \exp \{k_{16} \|\mathbf{b}\| + k_{17} \|\mathbf{Y}\| + k_{18}\} \times B_{s0} \, d\mathbf{b} \right| \\
& \leq e^{k_{19} + k_{20} \|\mathbf{Y}\|},
\end{aligned}$$

and  $Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$

$$\begin{aligned}
& = \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_2(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
& \quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T \right\} d\Lambda_s(t) \right\} d\mathbf{b}, \\
& \geq \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} \mathbf{b}^T \mathbf{b} - \sum_{j=1}^{n_N} \frac{B_2(\boldsymbol{\beta}; \mathbf{b})}{A(D(t_j; \phi))} - \exp \{k_{21} \|\mathbf{b}\| + k_{22} \|\mathbf{Y}\| + k_{23}\} \times B_{s0} \right\} d\mathbf{b}, \\
& \geq k_{24} > 0.
\end{aligned}$$

Moreover,

$$\begin{aligned}
& \|\nabla_{\boldsymbol{\theta}} Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \\
& = \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left( \nabla_{\boldsymbol{\theta}} \frac{Q_2}{Q_3} \right) \right\| + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left( \frac{d}{dv} \left( \frac{Q_2}{Q_3} \right) \right) \right\| \\
& = \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left[ \left( \nabla_{\boldsymbol{\theta}} Q_2 \right) \frac{1}{Q_3} + Q_2 \frac{(-1)}{Q_3^2} (\nabla_{\boldsymbol{\theta}} Q_3) \right] \right\| \\
& \quad + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left[ \left( \frac{d}{dv} Q_2 \right) \frac{1}{Q_3} + Q_2 \frac{(-1)}{Q_3^2} \left( \frac{d}{dv} Q_3 \right) \right] \right\| \\
& = \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + \left( \nabla_{\boldsymbol{\theta}} Q_2 \right) \frac{Q_1}{Q_3} - \left( \nabla_{\boldsymbol{\theta}} Q_3 \right) \frac{Q_2}{Q_3^2} \right\| + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + \left( \frac{d}{dv} Q_2 \right) \frac{Q_1}{Q_3} - \left( \frac{d}{dv} Q_3 \right) \frac{Q_1 Q_2}{Q_3^2} \right\| \\
& \leq \left( \|\nabla_{\boldsymbol{\theta}} Q_1\| + \left| \frac{d}{dv} Q_1 \right| \right) \left| \frac{Q_2}{Q_3} \right| + \left( \|\nabla_{\boldsymbol{\theta}} Q_2\| + \left| \frac{d}{dv} Q_2 \right| \right) \left| \frac{Q_1}{Q_3} \right| + \left( \|\nabla_{\boldsymbol{\theta}} Q_3\| + \left| \frac{d}{dv} Q_3 \right| \right) \left| \frac{Q_1 Q_2}{Q_3^2} \right| \\
& \leq e^{k_{25} + k_{26} \|\mathbf{Y}\|},
\end{aligned}$$

for some positive constants  $k_{25}$  and  $k_{26}$ . Therefore, by the mean-value theorem, we conclude that, for any  $(v_1, \boldsymbol{\theta}_1, \Lambda_{s1})$  and  $(v_2, \boldsymbol{\theta}_2, \Lambda_{s2})$  in  $[0, \tau] \times \Theta \times \mathcal{A}$ ,

$$\begin{aligned}
& |Q(v_1, \mathbf{O}; \boldsymbol{\theta}_1, \Lambda_{s1}) - Q(v_2, \mathbf{O}; \boldsymbol{\theta}_2, \Lambda_{s2})| \\
& \leq e^{k_{25} + k_{26} \|\mathbf{Y}\|} \left\{ \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| + |\Lambda_{s1}(V_s) - \Lambda_{s2}(V_s)| + \int_0^{V_s} |\Lambda_{s1}(t) - \Lambda_{s2}(t)| dt + |v_1 - v_2| \right\} \quad (3.32)
\end{aligned}$$

holds for some positive constants  $k_{25}$  and  $k_{26}$ .

Applying Theorem 2.7.5 (p159) in van der Vaart and Wellner (1996) to our situation, the entropy number for the class  $\mathcal{A}$  satisfies  $\log N_{[\cdot]}(\varepsilon, \mathcal{A}, L_2(P)) \leq K/\varepsilon$ , where  $K$  is a constant. Thus, we can find  $\exp\{K/\varepsilon\}$  brackets,  $\{[L_j, U_j]\}$ , to cover the class  $\mathcal{A}$  such that  $\|U_j - L_j\|_{L_2(P)} \leq \varepsilon$  for each pair of  $[L_j, U_j]$ . On the other hand, we can further find a partition of  $[0, \tau] \times \Theta$ , say  $I_1 \cup I_2 \cup \dots$ , such that the number of partitions is of the order  $(1/\varepsilon)^{d+1}$ , and, for any  $(v_1, \boldsymbol{\theta}_1)$  and  $(v_2, \boldsymbol{\theta}_2)$  in the same partition, their Euclidean distance is less than  $\varepsilon$ . Therefore, the partition  $\{I_1, I_2, \dots\} \times \{[L_j, U_j]\}$  bracket covers  $[0, \tau] \times \Theta \times \mathcal{A}$ , and the total number of the partitions is of order  $(1/\varepsilon)^{d+1} \exp\{1/\varepsilon\}$ . Hence, from (3.32), for any  $I_l$  and  $[L_j, U_j]$ , the set of the functions  $\{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : (v, \boldsymbol{\theta}) \in S_l, \Lambda_s \in \mathcal{A}, \Lambda_s \in [L_j, U_j]\}$  can be bracket covered by

$$\begin{aligned}
& \left[ Q(v_l, \mathbf{O}; \boldsymbol{\theta}_l, \Lambda_{sl}) - e^{k_{25} + k_{26} \|\mathbf{Y}\|} \left\{ \varepsilon + |U_j(V_s) - L_j(V_s)| + \int_0^{V_s} |U_j(t) - L_j(t)| dt \right\}, \right. \\
& \left. Q(v_l, \mathbf{O}; \boldsymbol{\theta}_l, \Lambda_{sl}) + e^{k_{25} + k_{26} \|\mathbf{Y}\|} \left\{ \varepsilon + |U_j(V_s) - L_j(V_s)| + \int_0^{V_s} |U_j(t) - L_j(t)| dt \right\} \right], \quad (3.33)
\end{aligned}$$

where  $(v_l, \boldsymbol{\theta}_l)$  is a fixed point in  $I_l$  and  $\Lambda_{sj}$  is a fixed function in  $[L_j, U_j]$ . Note that the  $L_2(P)$  distance between these two functions in the above bracket (3.33) is less than  $O(\varepsilon)$ . Therefore, we have

$$N_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_{L_2(P)}) \leq O\left(\left(\frac{1}{\varepsilon}\right)^{d+1} e^{1/\varepsilon}\right).$$

Furthermore,  $\mathcal{F}$  has an  $L_2(P)$ -integrable covering function, which is equal to  $O($

$e^{k_{25}+k_{26}\|\mathbf{Y}\|}$ ). From Theorem 2.5.6 (p130) in van der Vaart and Wellner (1996),  $\mathcal{F}$  is P-Donsker.

Additionally, in the above derivation, we also note that all the functions in  $\mathcal{F}$  are bounded from below by  $e^{-k_{27}-k_{28}\|\mathbf{Y}\|}$  for some positive constants  $k_{27}$  and  $k_{28}$ .

### 3.5.3.2 Proof of P-Donsker property of $\mathcal{G}$

Recall that we defined the class

$$\begin{aligned} \mathcal{G} = & \left\{ l_{\theta}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2], \right. \\ & \left. \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, (\mathbf{h}_1, h_2) \in \mathcal{H} \right\}, \end{aligned}$$

where  $(\mathbf{h}_1^{\beta}, \mathbf{h}_1^{\phi}, \mathbf{h}_1^b, \mathbf{h}_1^{\psi}, \mathbf{h}_1^{\gamma})$  denote the corresponding components of  $\mathbf{h}_1$  for the parameters  $(\boldsymbol{\beta}, \boldsymbol{\phi}, \text{Vec}(\boldsymbol{\Sigma}_b), \boldsymbol{\psi}, \boldsymbol{\gamma})$ , respectively. We can write that for  $(\mathbf{h}_1, h_2) \in \mathcal{H}$ ,

$$\begin{aligned} & l_{\theta}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] \\ = & \left[ \mu_1(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 - \int_0^{V_s} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 d\Lambda_s(t) \right] + \Delta h_2(V_s) - \int_0^{V_s} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) h_2(t) d\Lambda_s(t), \end{aligned}$$

where

$$\begin{aligned} & \mu_1(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 \\ = & \left\{ \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_2) d\mathbf{b} \right\}^{-1} \times \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_2) \times \left[ \frac{\mathbf{b}^T \boldsymbol{\Sigma}_b^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_b^{-1} \mathbf{b}}{2} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_b^{-1} \mathbf{D}_b) \right. \\ & + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta} + \tilde{\mathbf{X}}_j \mathbf{b}) - B(\boldsymbol{\beta}; \mathbf{b})}{(A(D(t_j; \phi)))^2} \right) A'(D(t_j; \phi)) h_1^{D(t_j; \phi)} + C'(Y_j; D(t_j; \phi)) h_1^{D(t_j; \phi)} \right\} \\ & \left. + \sum_{j=1}^{n_N} \left( \frac{Y_j \mathbf{X}_j}{A(D(t_j; \phi))} \mathbf{h}_1^{\beta} - B'(\boldsymbol{\beta}; \mathbf{b}) \mathbf{h}_1^{\beta} \right) + \Delta_s[(\tilde{\mathbf{Z}}(V_s) \circ \mathbf{h}_1^{\psi})^T \mathbf{b} + \mathbf{Z}(V_s) \mathbf{h}_1^{\gamma}] \right] d\mathbf{b}, \end{aligned}$$

$$\begin{aligned}
& \mu_2(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 \\
&= \left\{ \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b} \right\}^{-1} \\
& \quad \times \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \times \exp \left\{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma} \right\} \times \left[ (\tilde{\mathbf{Z}}(t) \circ \mathbf{h}_1^\psi)^T \mathbf{b} + \mathbf{Z}(t)\mathbf{h}_1^\gamma \right] d\mathbf{b}, \\
& \mu_3(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \left\{ \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b} \right\}^{-1} \times \int_{\mathbf{b}} G(\mathbf{b}, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \times \exp \left\{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}(t)\boldsymbol{\gamma} \right\} d\mathbf{b},
\end{aligned}$$

$\mathbf{D}_b$  is the symmetric matrix such that  $\text{Vec}(\mathbf{D}_b) = \mathbf{h}_1^b$ ,  $A'(D(t_j; \phi))$  and  $C'(Y_j; D(t_j; \phi))$  are the derivatives of  $A(D(t_j; \phi))$  and  $C(Y_j; D(t_j; \phi))$  with respect to  $D(t_j; \phi)$  respectively, and  $B'(\boldsymbol{\beta}; \mathbf{b})$  is the derivative of  $B(\boldsymbol{\beta}; \mathbf{b})$  with respect to  $\boldsymbol{\beta}$ .

For  $l = 1, 2, 3$ , we denote  $\nabla_{\boldsymbol{\theta}} \mu_l$  and  $\nabla_{\Lambda_s} \mu_l[\delta \Lambda_s]$  as the derivatives of  $\mu_l$  with respect to  $\boldsymbol{\theta}$  and  $\Lambda_s$  along the path  $\Lambda_s + \varepsilon \delta \Lambda_s$ . Then, using the similar arguments done in Section 3.5.3.1, it is verified that  $\nabla_{\Lambda_s} \mu_l[\delta \Lambda_s] = \int_0^t \mu_{l+3}(u, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\delta \Lambda_s(u)$  and there exist two positive constants  $q_1$  and  $q_2$  such that

$$\sum_l \{|\mu_l| + |\nabla_{\boldsymbol{\theta}} \mu_l|\} \leq e^{q_1 + q_2} \|\mathbf{Y}\|$$

By the mean value theorem, we have that, for any  $(\boldsymbol{\theta}, \Lambda_s, \mathbf{h}_1, h_2)$  and  $(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s, \tilde{\mathbf{h}}_1, \tilde{h}_2)$  in  $\Xi \times \mathcal{H}$ ,

$$\begin{aligned}
& l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2] \\
&= l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \mathbf{h}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[h_2] \\
& \quad + l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2] \\
&= [l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T] \mathbf{h}_1 + [l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)][h_2] \\
& \quad + l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)([h_2] - [\tilde{h}_2]) \\
&= (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \left[ \frac{d}{d\boldsymbol{\theta}} l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s) \right]_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \mathbf{h}_1 + \left[ \frac{d}{d\Lambda_s} l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) \right]_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*}^T [\Lambda_s - \tilde{\Lambda}_s] \mathbf{h}_1
\end{aligned}$$

$$\begin{aligned}
& + (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \left[ \frac{d}{d\boldsymbol{\theta}} l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right] [h_2] + \left[ \frac{d}{d\Lambda_s} l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right]^T [\Lambda_s - \tilde{\Lambda}_s] [h_2] \\
& + l_{\theta}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s) ([h_2] - [\tilde{h}_2]) \\
& = (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \nabla_{\boldsymbol{\theta}} \mu_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \mathbf{h}_1 - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\Lambda_s^*(t) \mathbf{h}_1 \\
& + \int_0^{V_s} \mu_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \tilde{\Lambda}_s)(t) \\
& - \int_0^{V_s} \int_0^t \mu_5(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \tilde{\Lambda}_s)(u) \mathbf{h}_1 d\Lambda_s^*(t) \\
& - \int_0^{V_s} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T (\Lambda_s - \tilde{\Lambda}_s) \mathbf{h}_1 dt \\
& - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
& - \int_0^{V_s} \int_0^t \mu_6(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \tilde{\Lambda}_s)(u) h_2(t) d\Lambda_s^*(t) \\
& - \int_0^{V_s} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T (\Lambda_s - \tilde{\Lambda}_s) h_2(t) dt \\
& + \mu_1(\mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) - \int_0^{V_s} \mu_2(t, \mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) d\tilde{\Lambda}_s(t) \\
& + \Delta_s(h_2(V_s) - \tilde{h}_2(V_s)) - \int_0^{V_s} \mu_3(t, \mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s) (h_2(V_s) - \tilde{h}_2(V_s)) d\tilde{\Lambda}_s(t), \tag{3.34}
\end{aligned}$$

where  $(\boldsymbol{\theta}^*, \Lambda_s^*)$  is equal to  $\varepsilon^*(\boldsymbol{\theta}, \Lambda_s) + (1 - \varepsilon^*)(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)$  for some  $\varepsilon^* \in [0, 1]$ . Thus, we have

$$\begin{aligned}
& |l_{\theta}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\theta}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2]| \\
& \leq e^{q_1 + q_2} \|\mathbf{Y}\| \left\{ \|\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}\| + \|\mathbf{h}_1 - \tilde{\mathbf{h}}_1\| + |\Lambda_s(V_s) - \tilde{\Lambda}_s(V_s)| \right. \\
& \quad + \int_0^{\tau} |\Lambda_s(t) - \tilde{\Lambda}_s(t)| [dt + d|h_2(t)| + d|\tilde{h}_2(t)|] \\
& \quad \left. + |h_2(V_s) - \tilde{h}_2(V_s)| + \int_0^{\tau} |h_2(V_s) - \tilde{h}_2(V_s)| [\Lambda_s(t) - \tilde{\Lambda}_s(t)] \right\}, \tag{3.35}
\end{aligned}$$

where  $d|h_2(t)| = dh_2^+(t) + dh_2^-(t)$  and  $d|\tilde{h}_2(t)| = d\tilde{h}_2^+(t) + d\tilde{h}_2^-(t)$ . As done in Section 3.5.3.1, by applying Theorem 2.7.5 (p159) in van der Vaart and Wellner (1996), we note that for a set  $\mathcal{H}_2 = \{h_2 : \|h_2\|_V \leq B_1\}$ ,  $\log N_{[\cdot]}(\varepsilon, \mathcal{H}_2, L_2(P)) \leq K/\varepsilon$  for a constant  $B_1$  and any probability measure  $P$  where  $K$  is a constant. Thus, we can find  $\exp\{K/\varepsilon\}$  brackets,  $\{[L_j, U_j]\}$ , to cover the class  $\mathcal{H}_2$  such that  $\|U_j - L_j\|_{L_2(P)} \leq \varepsilon$  for each pair of

$[L_j, U_j]$ . On the other hand, we can further find a partition of  $\mathcal{H}_1 = \{\mathbf{h}_1 : \|\mathbf{h}_1\| \leq 1\}$ , say  $I_1 \cup I_2 \cup \dots$ , such that the number of partitions is of the order  $(1/\varepsilon)$ , and, for any  $\mathbf{h}_1$  and  $h_2$  in the same partition, their Euclidean distance is less than  $\varepsilon$ . Therefore, the partition  $\{I_1, I_2, \dots\} \times \{[L_j, U_j]\}$  bracket covers  $\mathcal{H}_1 \times \mathcal{H}_2$ , and the total number of the partitions is of order  $(1/\varepsilon) \exp\{1/\varepsilon\}$ . Then, we obtain

$$\log N_{[\cdot]}(\varepsilon, \mathcal{G}, L_2(P)) \leq O\left(\frac{1}{\varepsilon} + \log \varepsilon\right).$$

Moreover,  $\mathcal{G}$  has an  $L_2(P)$ -integrable covering function, which is equal to  $O(e^{q_1+q_2\|\mathbf{Y}\|})$ . Hence, from Theorem 2.5.6 (p130) in van der Vaart and Wellner (1996),  $\mathcal{G}$  is P-Donsker.

Additionally, from (3.35), we can calculate

$$\begin{aligned} & |l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2]| \\ & \leq e^{q_1+q_2\|\mathbf{Y}\|} \left\{ \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + |\Lambda_s(V_s) - \Lambda_{s0}(V_s)| + \int_0^\tau |\Lambda_s(t) - \Lambda_{s0}(t)| dt \right\} \\ & \quad + \left| \int_0^\tau \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d(\Lambda_s(t) - \Lambda_{s0}(t)) \right|. \end{aligned} \quad (3.36)$$

If  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \rightarrow 0$  and  $\sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \rightarrow 0$ , the above expression converges to zero uniformly. Therefore,

$$\sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \mathbf{P} \left[ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right]^2 \longrightarrow 0.$$

### 3.5.3.3 Derivative operator $\nabla U_{\xi_0}$

From (3.34) in the previous Section 3.5.3.2, we can obtain that

$$\begin{aligned} & l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \\ & = [l_\theta(\boldsymbol{\theta}, \Lambda_s)^T - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T] \mathbf{h}_1 + [l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})][h_2] \end{aligned}$$

$$\begin{aligned}
&= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} \mu_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \mathbf{h}_1 - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \mathbf{h}_1 d\Lambda_s^*(t) \\
&\quad + \int_0^{V_s} \mu_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \Lambda_{s0})(t) \\
&\quad - \int_0^{V_s} \int_0^t \mu_5(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \Lambda_{s0})(u) \mathbf{h}_1 d\Lambda_s^*(t) \\
&\quad - \int_0^{V_s} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \Lambda_{s0})(t) \\
&\quad - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
&\quad - \int_0^{V_s} \int_0^t \mu_6(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \Lambda_{s0})(u) h_2(t) d\Lambda_s^*(t) \\
&\quad - \int_0^{V_s} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T h_2(t) d(\Lambda_s - \Lambda_{s0})(t) \\
&= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \left\{ \nabla_{\boldsymbol{\theta}} \mu_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\Lambda_s^*(t) \right\} \mathbf{h}_1 \\
&\quad + \mathbf{h}_1^T \left\{ \int_0^\tau I(t \leq V_s) [\mu_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \mu_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \right. \\
&\quad \quad \left. - \mu_5(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \int_t^{V_s} d\Lambda_s^*(u)] d(\Lambda_s - \Lambda_{s0})(t) \right\} \\
&\quad - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^\tau I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
&\quad - \int_0^\tau \left\{ I(t \leq V_s) \mu_6(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \int_0^{V_s} h_2(u) d\Lambda_s^*(u) \right. \\
&\quad \quad \left. + I(t \leq V_s) \mu_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) \right\} d(\Lambda_s - \Lambda_{s0})(t). \tag{3.37}
\end{aligned}$$

Then, we have that

$$\begin{aligned}
&\nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\
&= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{P} \left\{ \nabla_{\boldsymbol{\theta}} \mu_1(\mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\Lambda_{s0}(t) \right\} \mathbf{h}_1 \\
&\quad + \mathbf{h}_1^T \left\{ \int_0^\tau \mathbf{P} \left[ I(t \leq V_s) \left( \mu_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \mu_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right. \right. \right. \\
&\quad \quad \left. \left. - \mu_5(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u) \right) \right] d(\Lambda_s - \Lambda_{s0})(t) \right\} \\
&\quad - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^\tau \mathbf{P} \{ I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \} h_2(t) d\Lambda_{s0}(t)
\end{aligned}$$

$$\begin{aligned}
& - \int_0^\tau \mathbf{P} \left\{ I(t \leq V_s) \mu_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_0^{V_s} h_2(u) d\Lambda_{s0}(u) \right. \\
& \quad \left. + I(t \leq V_s) \mu_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) h_2(t) \right\} d(\Lambda_s - \Lambda_{s0})(t).
\end{aligned}$$

By the similar algebra done in (3.36), we can verify that, for  $j = 1, \dots, 6$ ,

$$\sup_{t \in [0, \tau]} \|\mu_j(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \mu_j(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\| \leq e^{q_3 + q + 4\|\mathbf{Y}\|} \left\{ \|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s^* - \Lambda_{s0}| \right\},$$

which implies that the linear operator  $\nabla U_{\xi_0}$  is bounded.

Then, we obtain that

$$\begin{aligned}
& \mathbf{P} [l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2]] \\
& = \nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] + o(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s - \Lambda_{s0}|)(\|\mathbf{h}_1\| + \|h_2\|_V).
\end{aligned}$$

Therefore,  $U_\xi$  is *Fréchet* differentiable at  $\xi_0$ .

Additionally, from (3.37) and the above expression, we have

$$\nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\Lambda_s - \Lambda_{s0})(t),$$

where

$$\begin{aligned}
\Omega_1[\mathbf{h}_1, h_2] &= \mathbf{P} \left\{ \nabla_{\boldsymbol{\theta}} \mu_1(\mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \mu_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\Lambda_{s0}(t) \right\} \mathbf{h}_1 \\
&\quad - \int_0^\tau \mathbf{P} \left\{ I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \mu_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} h_2(t) d\Lambda_{s0}(t)
\end{aligned}$$

and



$$\begin{aligned}
& \Omega_2[\mathbf{h}_1, h_2] \\
&= \mathbf{h}_1^T \mathbf{P} \left\{ I(t \leq V_s) [\mu_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \mu_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \mu_5(u, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u)] \right\} \\
&\quad - \mathbf{P} \left\{ I(t \leq V_s) \mu_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_0^{V_s} h_2(u) d\Lambda_{s0}(u) \right\} \\
&\quad - \mathbf{P} \left\{ I(t \leq V_s) \mu_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} h_2(t).
\end{aligned}$$

Thus,  $\Omega = (\Omega_1, \Omega_2)$  is the bounded linear operator from  $R^d \times BV[0, \tau]$  to itself. Furthermore, we note that  $\Omega = \mathbf{H} + (\mathbf{M}_1, \mathbf{M}_2)$ , where

$$\begin{aligned}
\mathbf{H}(\mathbf{h}_1, h_2) &= (\mathbf{h}_1, -\mathbf{P} \{ I(t \leq V_s) \mu_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \} h_2(t)), \\
\mathbf{M}_1(\mathbf{h}_1, h_2) &= \Omega_1[\mathbf{h}_1, h_2] - \mathbf{h}_1, \\
\mathbf{M}_2(\mathbf{h}_1, h_2) &= \mathbf{h}_1^T \mathbf{P} \left\{ I(t \leq V_s) [\mu_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \mu_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right. \\
&\quad \left. - \mu_5(u, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u)] \right\} \\
&\quad - \mathbf{P} \left\{ I(t \leq V_s) \mu_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_0^{V_s} h_2(u) d\Lambda_{s0}(u) \right\},
\end{aligned}$$

and also note that  $\mathbf{H}$  is obviously invertible. Since  $\mathbf{M}_1$  maps into a finite-dimensional space, it is compact. The image of  $\mathbf{M}_2$  is a continuously differentiable function in  $[0, \tau]$ . By the Arzela-Ascoli theorem (p41) in van der Vaart and Wellner (1996),  $\mathbf{M}_2$  is a compact operator from  $R^d \times BV[0, \tau]$  to  $BV[0, \tau]$ . Thus, we conclude that  $\Omega$  is the summation of an invertible operator  $\mathbf{H}$  and a compact operator  $\mathbf{M} = (\mathbf{M}_1, \mathbf{M}_2)$ .

### 3.6 Simulation Studies

In this section, we present some results from our simulation studies. Two sets of simulations with different generalized linear mixed models for the longitudinal outcomes

are performed. Binary and Poisson data are considered for longitudinal process in the first and second sets of simulations, respectively.

### 3.6.1 Binary longitudinal outcomes and survival time

In this first set of simulations, we assume  $Y_{ij}$  to be a binary outcome following

$$P(Y_{ij} = y_{ij}|b_i) = \exp \left\{ y_{ij}\eta_{ij} - \log(1 + \exp\{\eta_{ij}\}) \right\}, \quad y_{ij} = 0, 1,$$

with  $\eta_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + b_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3ij} + b_i$  for  $j = 1, \dots, n_i$ , and

$$h(t|b_i) = \lambda(t) \exp\{\psi b_i + \mathbf{Z}_i(t)\boldsymbol{\gamma}\} = \lambda(t) \exp\{\psi b_i + \gamma_1 Z_{1i} + \gamma_2 Z_{2i}\},$$

where  $b_i \sim N(0, \sigma_b^2)$ ,  $X_{1i} \equiv Z_{1i}$  are simulated from a Bernoulli distribution with success probability being 0.5, and  $X_{2i} \equiv Z_{2i}$  are simulated from the uniform distribution between 0 and 1. The longitudinal data are generated for every 0.3 unit of time, and thus  $X_{3ij}$ , the time at measurement, has the value of every 0.3 unit ranging over 0 through 2.4. We consider different  $\psi$  values of -0.1, 0, and 0.1 for negative, zero, and positive dependency between longitudinal process and survival time model, respectively. The parameters in the two models are chosen as  $\beta_0 = -1$ ,  $\beta_1 = 1$ ,  $\beta_2 = -0.5$ ,  $\beta_3 = -0.2$ ,  $\sigma_b^2 = 0.5$ ,  $\psi = -0.1/0/0.1$ ,  $\gamma_1 = -0.1$ ,  $\gamma_2 = 0.1$ , and  $\lambda(t) = 1$ . Censoring time is generated from the uniform distribution between 0.4 and 2.4, and the censoring proportion is around 25~35%. We consider different sample sizes ( $n=200, 400$ ) with 1000 replications. The average number of longitudinal observations ( $n_i$ ) is 3 with the range of 1 to 8. For the comparison of the estimated baseline cumulative hazards over simulations, we consider the three time points of 0.9, 1.4, 1.9 which are three quartiles of the observed survival time. The results of the maximum likelihood estimates for  $\boldsymbol{\theta}$  and baseline cumulative hazards at the given three time points and their respective standard error estimates

are reported in Table 3.3. The simulation study is conducted using R.

In Table 3.3, “True” gives the true values of parameters; the averages of the maximum likelihood estimates from the EM algorithm are in “Est.”; the sample standard deviations from 1000 simulations are reported in “SSD”; “ESE” is the average of 1000 standard error estimates based on the observed information matrix; “CP” is the coverage proportion of 95% nominal confidence intervals based on the estimated standard error “ESE”. Satterthwaite method is used for the coverage probability of  $\sigma_b^2$ .

From Table 3.3, we can see that even for the smaller sample size ( $n=200$ ), the bias of the estimates from EM algorithm is negligible for most cases. The estimated standard errors calculated from the observed information matrix are close to the sample standard deviations from the 1000 estimates, and the 95% confidence interval coverage rates are close to 0.95 except those for  $\psi$ . The parameter  $\psi$  tends to be underestimated with higher than the nominal level coverage rates, but the coverage rate is improved for larger sample size. Thus, with small sample size, the test for  $\psi$  is conservative, which strengthens the test results when rejecting the null ( $\psi = 0$ ), and the type I error becomes closer to the nominal level as sample size increases. In addition, the simulations show that the variances of the estimators decrease as the sample size ( $n$ ) increases. We also can see that the estimates are fairly robust and close to the true values for all different  $\psi$  values.

### 3.6.2 Poisson longitudinal outcomes and survival time

In the second set of simulations, we assume  $Y_{ij}$  to follow a Poisson distribution,

$$P(Y_{ij} = y_{ij}|b_i) = \exp \{y_{ij}\eta_{ij} - \exp\{\eta_{ij}\} - \log(y_{ij}!)\},$$

Table 3.3: Summary of simulation results of maximum likelihood estimation for binary longitudinal outcomes and survival time.

			n=200				n=400			
$\psi$	Par.	True	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP
- .1	$\beta_0$	-1.0	-1.008	.270	.272	.951	-1.006	.194	.191	.945
	$\beta_1$	1.0	1.015	.241	.232	.942	.997	.162	.162	.949
	$\beta_2$	- .5	- .522	.405	.392	.939	- .500	.282	.274	.945
	$\beta_3$	- .2	- .173	.252	.245	.947	- .187	.181	.172	.952
	$\sigma_b^2$	.5	.502	.231	.288	.968	.493	.168	.200	.974
	$\psi$	- .1	- .083	.353	.390	.990	- .102	.235	.245	.978
	$\gamma_1$	- .1	- .102	.174	.174	.949	- .104	.121	.121	.955
	$\gamma_2$	.1	.099	.296	.303	.953	.097	.214	.210	.950
	$\Lambda(.9)$	.9	.910	.181	.184	.955	.906	.132	.128	.943
	$\Lambda(1.4)$	1.4	1.444	.294	.299	.962	1.421	.211	.204	.943
	$\Lambda(1.9)$	1.9	1.980	.446	.449	.955	1.945	.312	.302	.951
0	$\beta_0$	-1.0	-1.012	.281	.273	.944	-1.008	.192	.192	.948
	$\beta_1$	1.0	1.014	.235	.233	.954	1.006	.164	.163	.948
	$\beta_2$	- .5	- .501	.414	.393	.934	- .503	.278	.276	.949
	$\beta_3$	- .2	- .190	.263	.246	.936	- .194	.175	.173	.950
	$\sigma_b^2$	.5	.505	.237	.290	.960	.503	.176	.203	.963
	$\psi$	.0	.008	.375	.388	.994	.003	.236	.241	.980
	$\gamma_1$	- .1	- .108	.180	.174	.942	- .103	.113	.121	.968
	$\gamma_2$	.1	.098	.309	.303	.949	.101	.209	.210	.951
	$\Lambda(.9)$	.9	.920	.188	.186	.952	.905	.131	.127	.944
	$\Lambda(1.4)$	1.4	1.463	.306	.303	.948	1.415	.206	.202	.952
	$\Lambda(1.9)$	1.9	2.006	.462	.457	.953	1.937	.306	.299	.948
.1	$\beta_0$	-1.0	-1.009	.285	.274	.948	-1.000	.192	.192	.945
	$\beta_1$	1.0	1.004	.224	.234	.964	1.004	.166	.163	.952
	$\beta_2$	- .5	- .510	.414	.395	.945	- .512	.284	.276	.943
	$\beta_3$	- .2	- .186	.260	.249	.948	- .185	.189	.175	.929
	$\sigma_b^2$	.5	.519	.249	.295	.946	.498	.174	.203	.965
	$\psi$	.1	.117	.354	.386	.990	.129	.246	.247	.986
	$\gamma_1$	- .1	- .096	.175	.175	.946	- .101	.117	.122	.966
	$\gamma_2$	.1	.086	.312	.304	.944	.101	.212	.211	.948
	$\Lambda(.9)$	.9	.915	.184	.185	.957	.904	.129	.127	.946
	$\Lambda(1.4)$	1.4	1.455	.305	.303	.955	1.413	.203	.203	.954
	$\Lambda(1.9)$	1.9	2.010	.481	.463	.952	1.938	.303	.302	.959

with  $\eta_{ij}$  defined as in Section 3.6.1. We also consider the same hazards model and simulation setting as those used in Section 3.6.1 except  $\sigma_b^2 = 0.2$ . The simulated Poisson longitudinal outcomes range over 0 to 7 with the average 0.5.

Table 3.4 shows that overall the estimates perform well even for the smaller sample size  $n = 200$  with small biases of the estimates except  $\psi$ . We conducted additional simulations with sample sizes of 800 and 1000, and the bias of  $\psi$  existing for the small sample size decreases as sample size increases over 200, 400, 800 and 1000. The estimated standard errors using the observed information matrix are close to the sample standard deviations, and the 95% confidence interval coverage rates are close to 0.95 except for  $\sigma_b^2$  and  $\psi$ .

From Table 3.4,  $\psi$  is seemingly underestimated with higher than the nominal coverage rates, but the coverage rate decreases to close to 95% nominal level as sample size increases. Additional simulations we conducted show that, with sample sizes of 800, the 95% confidence interval coverage rates for  $\psi = -0.1, 0$  and  $0.1$  were 95.5%, 95.9% and 95.9%, respectively.  $\sigma_b^2$  also appears to have high coverage rates, which may be due to numerical problem since its coverage rates are still high for larger sample sizes. This implies that variance of  $\sigma_b^2$  may not be estimated well for Poisson longitudinal distribution. In the meantime, the test for  $\sigma_b^2$  is conservative, which strengthens the test result for rejecting the null ( $\sigma_b^2 = 0$ ). On the other hand, profile likelihood may be an alternative estimation approach for  $\sigma_b^2$ . It is also shown that the variances of the estimators decrease for larger sample size, and the estimates are fairly robust and close to the true values for all three different  $\psi$  values.

### 3.7 Analysis of the CHANCE Study

We now return to the CHANCE study described in Section 3.2, and apply our proposed method to Head and Neck Cancer Specific symptoms (HNCS) among QoL domains with

Table 3.4: Summary of simulation results of maximum likelihood estimation for Poisson longitudinal outcomes and survival time.

			n=200				n=400			
$\psi$	Par.	True	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP
- .1	$\beta_0$	-1.0	-1.001	.186	.192	.959	-1.005	.135	.135	.949
	$\beta_1$	1.0	1.014	.165	.161	.949	1.008	.118	.113	.933
	$\beta_2$	- .5	- .513	.275	.260	.938	- .502	.189	.182	.951
	$\beta_3$	- .2	- .188	.178	.164	.941	- .189	.128	.115	.921
	$\sigma_b^2$	.2	.195	.074	.096	.978	.197	.051	.067	.986
	$\psi$	- .1	- .056	.603	.621	.981	- .075	.417	.388	.970
	$\gamma_1$	- .1	- .098	.177	.175	.952	- .098	.126	.121	.940
	$\gamma_2$	.1	.081	.311	.305	.946	.103	.218	.211	.956
	$\Lambda(.9)$	.9	.922	.189	.187	.949	.902	.130	.127	.946
	$\Lambda(1.4)$	1.4	1.466	.321	.307	.941	1.418	.209	.204	.950
	$\Lambda(1.9)$	1.9	2.036	.502	.474	.950	1.947	.308	.304	.950
0	$\beta_0$	-1.0	-1.000	.190	.192	.946	-1.003	.134	.135	.948
	$\beta_1$	1.0	1.009	.162	.161	.944	1.004	.114	.113	.949
	$\beta_2$	- .5	- .512	.277	.260	.933	- .496	.187	.183	.940
	$\beta_3$	- .2	- .190	.179	.165	.934	- .187	.125	.116	.935
	$\sigma_b^2$	.2	.195	.076	.096	.977	.199	.052	.067	.984
	$\psi$	.0	.017	.606	.628	.989	.061	.403	.382	.968
	$\gamma_1$	- .1	- .098	.174	.175	.952	- .105	.125	.121	.940
	$\gamma_2$	.1	.094	.311	.305	.952	.101	.212	.211	.950
	$\Lambda(.9)$	.9	.914	.185	.185	.955	.906	.130	.128	.944
	$\Lambda(1.4)$	1.4	1.459	.305	.304	.954	1.428	.209	.205	.942
	$\Lambda(1.9)$	1.9	2.016	.472	.464	.957	1.961	.308	.305	.945
.1	$\beta_0$	-1.0	-1.000	.192	.193	.947	-1.005	.131	.135	.958
	$\beta_1$	1.0	1.013	.168	.162	.939	1.007	.115	.114	.945
	$\beta_2$	- .5	- .522	.279	.261	.937	- .500	.191	.183	.939
	$\beta_3$	- .2	- .193	.191	.167	.929	- .188	.129	.117	.928
	$\sigma_b^2$	.2	.198	.076	.096	.978	.197	.053	.067	.984
	$\psi$	.1	.152	.609	.627	.993	.137	.403	.390	.975
	$\gamma_1$	- .1	- .098	.174	.176	.944	- .103	.120	.121	.951
	$\gamma_2$	.1	.089	.303	.306	.953	.090	.203	.211	.958
	$\Lambda(.9)$	.9	.919	.187	.187	.945	.913	.122	.128	.962
	$\Lambda(1.4)$	1.4	1.461	.317	.308	.952	1.439	.205	.207	.961
	$\Lambda(1.9)$	1.9	2.038	.519	.480	.951	1.965	.308	.307	.956

survival time in this section. We are interested in testing the correlation between survival time and longitudinal QoL outcomes and investigating which variables predict the QoL satisfaction or the risk of death or both. In the full models for both longitudinal QoL and survival time, race, the number of 12 oz. beers consumed per week, household income, surgery, radiation therapy, chemotherapy, tumor site, and tumor stage are considered as categorical, and age at diagnosis, the number of persons supported by household income, body mass index (BMI), and the total number of medical conditions reported as continuous. Additionally, 2 interactions with race, i.e. race  $\times$  the total number of medical conditions reported and race  $\times$  tumor site, are included in both models since we are particularly interested in the difference of QoL and survival between African American and White. Time at survey measurement is also included as a covariate for longitudinal outcomes. A random intercept for the dependence between the QoL satisfaction and the risk of death is included in both models, and assumed to follow a normal distribution with mean zero. In addition to the simultaneous analysis, we also conduct separate analyses fitting the generalized linear mixed model and the Cox proportional hazards model to the longitudinal QoL and survival time respectively and compare the results to those from our proposed simultaneous method.

After fitting the simultaneous models with all the covariates, we use backward variable selection based on the Likelihood Ratio Test (LRT) and find that surgery, chemotherapy, tumor site, age at diagnosis, and all 2 interactions are not statistically significant in both models for HNCS QoL satisfaction and survival time at the significance level 0.05. We remove these variables and refit the simultaneous models. Then, the LRT shows that race, radiation therapy, the number of persons supported by household income, BMI, and the total number of medical conditions reported are not statistically significant for the risk of death. We further reduce the models by removing them from the hazards model and refit the reduced simultaneous models. Table 3.5

gives the results from this final models. From the “Simultaneous” columns, we see that the number of 12 oz. beers consumed per week, household income, tumor stage, and the total number of medical conditions reported are significantly associated with both patients’ HNCS QoL satisfaction and hazard of death. Using 30 or more of 12 oz. beers consumed per week as the reference group, all categories of the smaller amount are associated with HNCS QoL satisfaction and lower risk of death, higher household income is overall associated with HNCS QoL satisfaction and lower risk of death, and both patients’ HNCS QoL satisfaction and risk of death are significantly different for patients in different tumor stages. Specifically, for instance, with the log-scaled odds and hazard ratios of 1.060 and -1.076 for HNCS QoL satisfaction and death respectively, patients who consumed 5 to 14 of 12 oz. beers per week appear to have 2.886 times odds for HNCS QoL satisfaction and 0.341 times hazards of death compared to those that consumed 30 or more of 12 oz. beers per week in the study after adjusting for the other covariates in the model. Looking at the number of medical conditions reported, for each additional medial condition reported, the odds ratio of HNCS QoL satisfaction is decreased by 16% and the hazard of death is increased by 29%. That is, patients with a greater number of medical conditions reported have lower HNCS QoL satisfaction and higher risk of death after adjusting for the other covariates in the model. In the meantime, race (African-American), radiation therapy, the number of persons supported by household income, and BMI are selected only in the HNCS QoL longitudinal model. African-Americans, patients not treated with radiation therapy, patients in the family with the smaller number of persons supported by household income, or patients with higher BMI are likely to be more satisfied with longitudinal HNCS QoL while the risk of death is not affected by race, radiation therapy, the number of persons supported by household income and BMI. Furthermore, we also find that time at survey measurement is statistically significant in the HNCS QoL longitudinal model implying that



patients are more satisfied over time. The parameter  $\psi$  for the dependence between longitudinal HNCS QoL and survival time is negative and has p-value as 0.131. This implies that the longitudinal HNCS QoL and survival time are marginally correlated and some latent factors which increase HNCS QoL satisfaction also decrease the risk of death. For the purpose of comparison, we conducted separate analyses for longitudinal HNCS QoL and survival time whose results are given in the last three columns of Table 3.5. The generalized linear mixed model (GLMM) and the Cox proportional hazards model are used for longitudinal outcomes and survival time respectively. The GLMM also considers individual heterogeneity through subject-specific random effects although it does not incorporate the correlation between longitudinal outcomes and survival time. Comparing the results from the simultaneous and separate analyses of Table 3.5, we can see our simultaneous analysis additionally indicates the number of persons supported by household income, BMI, and the total number of medical conditions reported in the HNCS QoL longitudinal model (p-values=0.025, 0.007, and 0.030, respectively) and the number of 12 oz. beers consumed per week in the hazard model (p-values=0.045 and 0.005 for ‘None’ and ‘5 to 14’) as significant which are not selected by separate analyses. Figure 3.1 shows the estimated baseline cumulative hazard rates over follow-up time with the 95% confidence interval. Since the baseline cumulative hazard rates are bounded by 0, we first log-transformed the estimated baseline cumulative hazard rates and obtained the 95% lower and upper bounds for the log-scaled estimated baseline cumulative hazards. Then, we re-transformed them into their original scale. The estimated baseline cumulative hazard rates look flat at the very early time within a year, but soon appear to be linearly increasing. Figure 3.2 shows the Kaplan-Meier estimates (solid line) and the predicted survival probabilities based on the simultaneous analysis (dashed line). These two survival curves are very close to each other which implies our proposed method fits the data well.

Table 3.5: Analyses results for the HNCS QoL and survival time of the CHANCE study

Parameter		Simultaneous			Separate		
		Est.	ESE	P-value	Est.	ESE	P-value
< HNCS QoL longitudinal model >							
Intercept	$\beta_0$	.744	.538	.167	1.190	.390	.002
Race (ref= White)							
– African American	$\beta_1$	.564	.229	.014	.511	.256	.047
# of 12 oz. beers consumed per week (ref=30 or more)							
– None	$\beta_2$	.636	.269	.018	.622	.300	.038
– less than 1	$\beta_3$	.830	.357	.020	.735	.396	.064
– 1 to 4	$\beta_4$	1.302	.294	<.001	1.268	.326	<.001
– 5 to 14	$\beta_5$	1.060	.251	<.001	1.018	.279	<.001
– 15 to 29	$\beta_6$	.601	.289	.037	.547	.327	.095
Household income (ref= level1: 0–10K)							
– level2: 20–30K	$\beta_7$	-.271	.231	.241	-.328	.258	.204
– level3: 40–50K	$\beta_8$	.297	.255	.245	.250	.282	.376
– level4: ≥ 60K	$\beta_9$	1.199	.274	<.001	1.045	.286	<.001
Radiation therapy (ref= No)							
– Yes	$\beta_{10}$	-1.132	.260	<.001	-1.048	.280	<.001
Tumor stage (ref= I)							
– II	$\beta_{11}$	-.416	.300	.166	-.352	.330	.286
– III	$\beta_{12}$	-1.335	.284	<.001	-1.198	.314	<.001
– IV	$\beta_{13}$	-1.175	.254	<.001	-1.057	.277	<.001
# of persons supported by household income	$\beta_{14}$	-.189	.084	.025			
BMI	$\beta_{15}$	.041	.015	.007			
Total # of medical conditions reported	$\beta_{16}$	-.175	.080	.030			
Time at survey measurement (years)	$\beta_{17}$	.241	.066	<.001	.254	.067	<.001
variance of random effects	$\sigma_b^2$	.303	.173	.013	1.169	.257	
< Hazards model >							
Random effect coefficient	$\psi$	-1.427	.946	.131			
# of 12 oz. beers consumed per week (ref=3 or more)							
– None	$\gamma_1$	-.772	.386	.045			
– less than 1	$\gamma_2$	-.155	.426	.715			
– 1 to 4	$\gamma_3$	-.802	.414	.053			
– 5 to 14	$\gamma_4$	-1.076	.383	.005			
– 15 to 29	$\gamma_5$	-.591	.399	.139			
Household income (ref= level1: 0–10K)							
– level2: 20–30K	$\gamma_6$	-.218	.294	.459	-.219	.263	.406
– level3: 40–50K	$\gamma_7$	-.941	.371	.011	-.928	.331	.005
– level4: ≥ 60K	$\gamma_8$	-1.463	.401	<.001	-1.393	.358	<.001
Tumor stage (ref= I)							
– II	$\gamma_9$	-.199	.465	.668	-.295	.435	.498
– III	$\gamma_{10}$	.235	.433	.588	.136	.389	.727
– IV	$\gamma_{11}$	1.059	.360	.003	.914	.295	.002
Total # of medical conditions reported	$\gamma_{12}$	.256	.110	.020	.205	.091	.025

P-value for testing  $\sigma_b^2$  being zero is based on a mixture of 0 and  $\chi^2$  distribution with 1 degree of freedom with equal mixing probabilities.

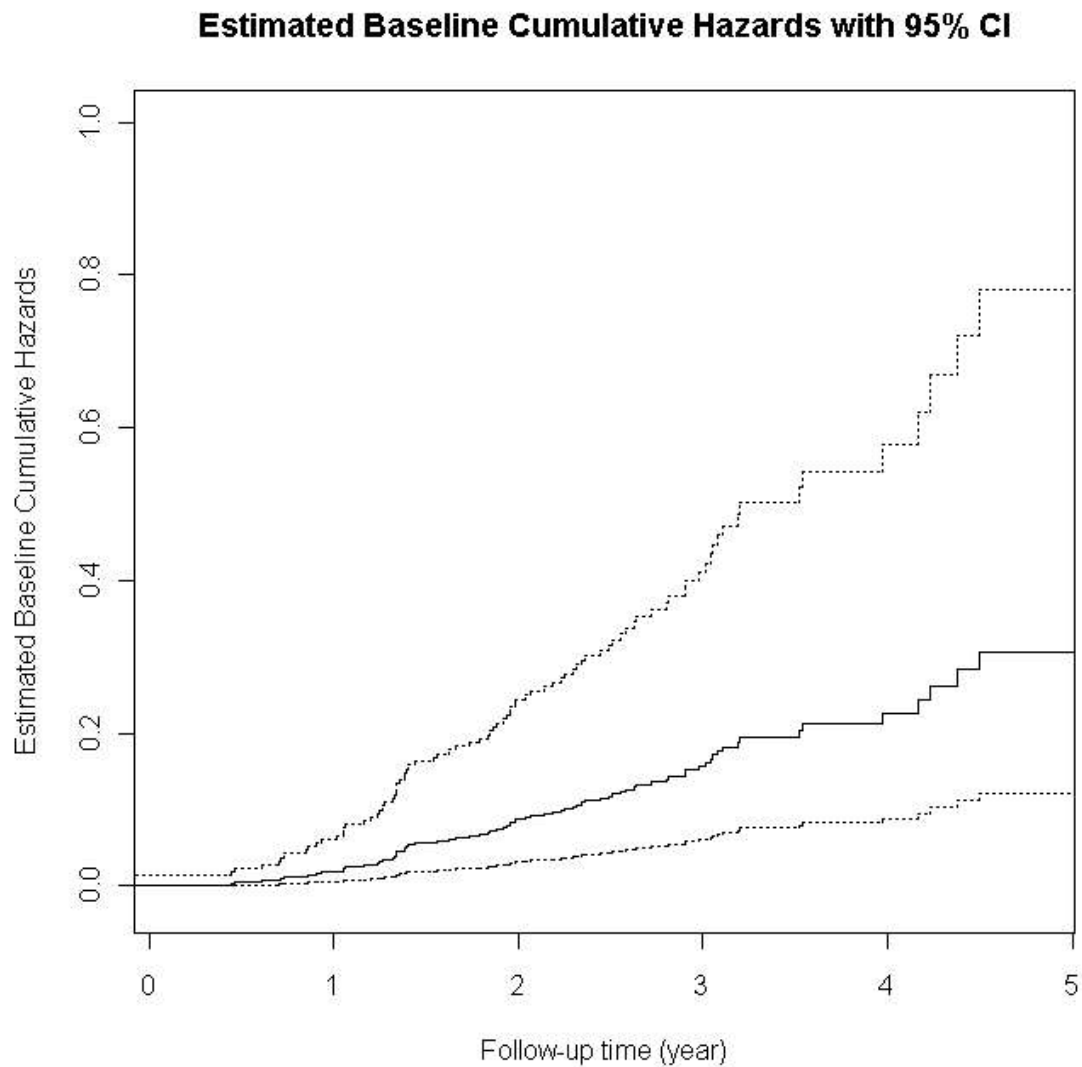


Figure 3.1: Estimated baseline cumulative hazards (solid line) with 95% confidence interval (dashed lines) by the simultaneous analysis of HNCS QoL longitudinal outcome and survival time

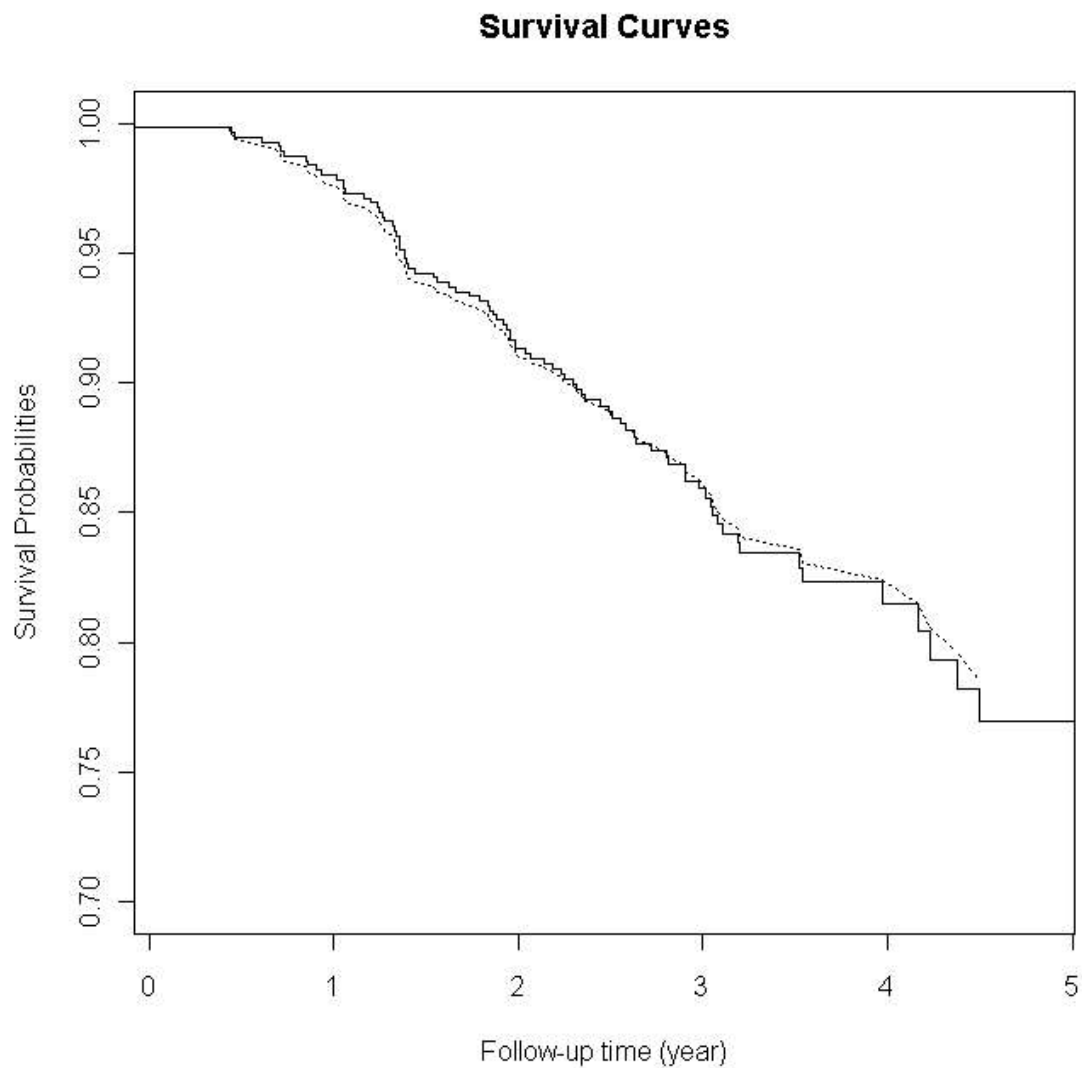


Figure 3.2: Kaplan-Meier estimates (solid line) and the predicted survival probabilities based on the simultaneous analysis of HNCS QoL longitudinal outcome and survival time (dashed line)

### 3.8 Concluding Remarks

We have proposed a method for the simultaneous modeling of longitudinal outcomes including both categorical and continuous data with a generalized linear mixed model and survival time with a stratified multiplicative proportional hazards model through random effects. We have also developed a maximum likelihood estimation method for the proposed simultaneous model, and presented asymptotic properties of the proposed estimators. The proposed estimation procedure using EM algorithm has been assessed via simulation studies. The proposed estimates performed well in finite samples. The variance estimates based on the observed information matrix approximate the true variance well in finite samples.

The proposed method was applied to the CHANCE study data. The results for longitudinal HNCS and survival time have shown that, after adjusting for the other covariates in the simultaneous models, the lower amount of beers consumed per week, higher household income, lower stage, and the lower total number of medical conditions reported are associated with more HNCS QoL satisfaction and lower risk of death. Further, African-Americans, patients not treated with radiation therapy, patients in the family with the smaller number of persons supported by household income, or patients with higher BMI are likely to be more satisfied with longitudinal HNCS QoL while the risk of death is not affected by race, radiation therapy, the number of persons supported by household income and BMI. Time at survey measurement in the HNCS QoL longitudinal model is also statistically significant implying that patients are more satisfied over time. Furthermore, our proposed method additionally finds more predictors including: the number of persons supported by household income, BMI, and the total number of medical conditions reported in the HNCS QoL longitudinal model and the predictor, number of 12 oz. beers consumed per week, in the hazard model while separate analyses do not select them. This result comparing the simultaneous and separate

analyses supports that, even when the longitudinal outcomes and survival time are only marginally correlated, our simultaneous analysis could provide better power than separate analyses not considering the dependency between the longitudinal outcomes and survival time.

In our proposed method, all the information on survival, longitudinal outcomes, and covariates are used. As a result of this, the parameter estimates can be more efficient. The proposed model also generalizes previous work to general longitudinal outcomes. This work fills in some gaps in the joint modeling research. Future work can include relaxing normal assumption for the random effects and considering generalization to mixed types of longitudinal outcomes.

# Chapter 4

## JOINT MODELING OF SURVIVAL TIME AND LONGITUDINAL OUTCOMES WITH FLEXIBLE RANDOM EFFECTS

### 4.1 Introduction

In biomedical or public health research, it is common that both longitudinal outcomes over time and survival endpoint are collected for the same subject along with the subject's characteristics or risk factors. Investigators are interested in finding important variables which predict both longitudinal outcomes and survival time. Among the existing approaches for longitudinal data and survival time, the selection model and the pattern mixture model have been widely used. The selection model estimating the distribution of survival time given longitudinal data was studied by numerous authors, for example, Tsiatis *et al.* (1995), Tsiatis and Davidian (2001), Xu and Zeger (2001a,b) and Tseng *et al.* (2005). The pattern mixture model focuses on the trend of longitudinal outcomes conditional on survival time and was studied by Wu and Carroll (1988), Hogan and Laird (1997), Albert and Follmann (2000, 2007) and Ding and Wang (2008) among others. On the other hand, simultaneous modeling of the longitudinal and survival data

was proposed by Xu and Zeger (2001b), Zeng and Cai (2005), Elashoff *et al.* (2007, 2008) and Rizopoulos *et al.* (2008). Wang and Taylor (2001), Brown and Ibrahim (2003) and Hu *et al.* (2009) studied simultaneous modeling in the Bayesian perspective.

In all the joint models, random effects are incorporated to accommodate the latent dependence between survival time and longitudinal outcomes. Random effects are conventionally assumed to be normally distributed. However, it is unclear whether the normality assumption is truly satisfied in practice. Furthermore, misspecifying normality assumption can lead to serious bias in estimation (Neuhaus *et al.*, 1992; Kleinman and Ibrahim, 1998; Heagerty and Kurland, 2001; Agresti *et al.*, 2004).

In this paper, we assume that the underlying distribution of random effects is unknown. In estimating model parameters, we propose to use a mixture of Gaussian distributions as an approximation for the unknown random effect distribution. Moreover, we simultaneously model the survival time with a stratified Cox proportional hazards model and longitudinal outcomes with a generalized linear mixed model to incorporate both categorical and continuous longitudinal outcomes. Finite sample properties of the proposed estimators and robustness of the mixture distribution are assessed via simulations. We adopt AIC and BIC for selecting the number of mixtures and also conduct simulation studies to assess these selection procedures.

The outline of this paper is as follows. In Section 4.2, we present a simultaneous modeling for longitudinal outcomes and survival time with random effects from an unknown distribution, and describe the inference procedure. Asymptotic properties of the proposed estimators and the technical details of their proofs are investigated in Section 4.3 and Section 4.4, respectively. Numerical results from simulation studies are given in Section 4.5. Our proposed method is illustrated with the data from the Carolina Head and Neck Cancer Study (CHANCE) in Section 4.6. In Section 4.7, we discuss some further consideration. EM-algorithms are provided in Appendix.



## 4.2 Models and Inference Procedure

### 4.2.1 Model formulation and notation

We use  $Y(t)$  to denote the value of a longitudinal marker process at time  $t$ . Suppose  $Y(t)$  is from a distribution belonging to exponential family in order to incorporate both continuous and categorical measurements. Let  $T$  denote survival time, and suppose that the survival time  $T$  is possibly right censored. Suppose a set of  $n$  subjects are followed over an interval  $[0, \tau]$ , where  $\tau$  is the study end time. Denote  $\mathbf{b}_i^*$ ,  $i = 1, \dots, n$ , as a vector of subject-specific random effects of dimension  $d_b$  and  $\mathbf{b}_i^*$ 's are mutually independent.

Given the random effects  $\mathbf{b}_i^*$ , the observed covariates, and the observed outcome history till time  $t$ , we assume that the longitudinal outcome  $Y_i(t)$  at time  $t$  for subject  $i$  follows a distribution from the exponential family with density,

$$\exp \left\{ \frac{y_i \eta_i(t) - B(\eta_i(t))}{A(D_i(t; \phi))} + C(y_i, D_i(t; \phi)) \right\} \quad (4.1)$$

with  $\mu_i(t) = E(Y_i(t)|\mathbf{b}_i^*) = B'(\eta_i(t))$  and  $v_i(t) = \text{Var}(Y_i(t)|\mathbf{b}_i^*) = B''(\eta_i(t))A(D_i(t; \phi))$ , satisfying

$$\eta_i(t) = g(\mu_i(t)) = \mathbf{X}_i(t)\boldsymbol{\beta} + \tilde{\mathbf{X}}_i(t)\mathbf{b}_i^*$$

and  $v_i(t) = v(\mu_i(t))A(D_i(t; \phi))$ , where  $g(\cdot)$  and  $v(\cdot)$  are known link and variance functions respectively,  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  are the row vectors of the observed covariates for subject  $i$ , and  $\boldsymbol{\beta}$  is a column vector of coefficients for  $\mathbf{X}_i(t)$ .  $\mathbf{X}_i(t)$  does not include intercept and it does not contain any covariates in  $\tilde{\mathbf{X}}_i(t)$  because the intercept and any potential common covariates for fixed effects are combined with the corresponding random effects in  $\tilde{\mathbf{X}}_i(t)$  so that mean of random effects does not have restriction.

Given the random effects  $\mathbf{b}_i^*$ , the observed covariates, and the observed survival history before time  $t$ , the conditional hazard rate function for the survival time  $T_i$  of

subject  $i$  is assumed to follow a stratified multiplicative hazards model,

$$\lambda_s(t) \exp\{\tilde{\mathbf{Z}}_i(t)(\boldsymbol{\psi} \circ \mathbf{b}_i^*) + \mathbf{Z}_i(t)\boldsymbol{\gamma}\}, \quad (4.2)$$

where  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  are the row vectors of the observed covariates and may share some components,  $\boldsymbol{\psi}$  is a vector of parameters of the coefficients for random effects,  $\lambda_s(t)$  is the  $s$ -th stratum baseline hazard rate function, and  $\boldsymbol{\gamma}$  is a column vector of coefficients for  $\mathbf{Z}_i(t)$ . Note that  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  do not include dummy variables for strata since baseline hazard rate is stratum-specific. Here, for any vectors  $\mathbf{a}_1$  and  $\mathbf{a}_2$  of the same dimension,  $\mathbf{a}_1 \circ \mathbf{a}_2$  denotes the component-wise product. In addition,  $\tilde{\mathbf{X}}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  have the same dimensions as  $\mathbf{b}_i^*$ 's. For the subject-specific random effects  $\mathbf{b}_i^*$ , we assume the underlying distribution of  $\mathbf{b}_i^*$  is unknown and denote its density as  $f(\mathbf{b}_i^*)$ .

Let  $n_i$  be the number of the observed longitudinal measurements for subject  $i$ , and assume that the distributions of  $n_i$  and the observation times for longitudinal measurements are independent of the parameters of interest in this joint model. The observed data from  $n$  subjects are  $(n_i, Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$ ,  $j = 1, \dots, n_i$ ,  $i = 1, \dots, n$ , and  $(V_i, \Delta_i, S_i, \{(\mathbf{Z}_i(t), \tilde{\mathbf{Z}}_i(t)) : t \leq V_i\})$ ,  $i = 1, \dots, n$ , where for subject  $i$ ,  $(Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$  is the  $j$ -th observation of  $(Y_i(t), \mathbf{X}_i(t), \tilde{\mathbf{X}}_i(t))$ ,  $C_i$  is the right-censoring time and independent of  $T_i$  and  $Y_i(t)$  given the covariates and the random effects,  $V_i = \min(T_i, C_i)$ ,  $S_i$  denotes the stratum, and  $\Delta_i = I(T_i \leq C_i)$ . For all  $n$  subjects, we write  $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_n^T)^T$ ,  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T$ ,  $\mathbf{V} = (V_1, \dots, V_n)^T$ , and  $\mathbf{b}^* = (\mathbf{b}_1^{*T}, \dots, \mathbf{b}_n^{*T})^T$ . Then, the likelihood function of the complete data  $(\mathbf{Y}, \mathbf{V}, \mathbf{b}^*)$  has the form,

$$\begin{aligned} L_c(\mathbf{Y}, \mathbf{V}, \mathbf{b}^*) &= \prod_{i=1}^n f(\mathbf{Y}_i | \mathbf{b}_i^*) \left( \prod_{s=1}^S [f(V_i | \mathbf{b}_i^*)]^{I(S_i=s)} \right) f(\mathbf{b}_i^*) \\ &= \prod_{i=1}^n \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i^*) - B(\boldsymbol{\beta}; \mathbf{b}_i^*)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \end{aligned}$$

$$\begin{aligned}
& \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i^*) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
& \quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i^*) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
& \times f(\mathbf{b}_i^*),
\end{aligned} \tag{4.3}$$

and the full likelihood function of the observed data  $(\mathbf{Y}, \mathbf{V})$  is expressed as

$$L_f(\mathbf{Y}, \mathbf{V}) = \int_{\mathbf{b}^*} L_c(\mathbf{Y}, \mathbf{V}, \mathbf{b}^*) d\mathbf{b}^*. \tag{4.4}$$

The parameter  $\boldsymbol{\psi}$  in model (4.2) characterizes the dependence between the longitudinal outcomes and the survival time due to latent random effects:  $\boldsymbol{\psi} = \mathbf{0}$  means the dependence between the survival time and longitudinal responses are not due to these latent variables;  $\boldsymbol{\psi} \neq \mathbf{0}$  means such dependence may be due to these latent variables. In other words,  $\boldsymbol{\psi} > \mathbf{0}$  implies there may be some latent factors increasing both the longitudinal outcomes and the risk of survival endpoint simultaneously while  $\boldsymbol{\psi} < \mathbf{0}$  implies some latent factors causing the increment of longitudinal outcomes may decrease the risk of survival endpoint.

### 4.2.2 Inference procedure

For parameter estimation, we approximate the random effect distribution with a mixture of Gaussian distributions. This method was studied in some literature to extend normality assumption of random effects. For instance, Verbeke and Lesaffre (1996), Verbeke and Molenberghs (2000), and Zhang and Davidian (2001) used it in a linear mixed model, and Verbeke and Lesaffre (1996), Fieuws *et al.* (2004), and Caffo *et al.* (2007) considered it in a GLMM. Alternatively, Ghidey *et al.* (2004) and Komárek and Lesaffre (2008a) used the penalized Gaussian mixture (PGM) approach in a linear mixed model and a GLMM respectively. Furthermore, Komárek and Lesaffre (2008b,

2009) suggested a Bayesian accelerated failure time (AFT) model with random effects following a PGM.

For the subject-specific random effects  $\mathbf{b}_i^*$  in Section 4.2.1, we approximate the distribution of  $\mathbf{b}_i^*$  with a mixture of a finite number of  $d_b$ -dimensional multivariate normal distributions. That is, the distribution of  $\mathbf{b}_i^*$  is approximated by  $\sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_b)$ , where  $K$  is the number of mixture components. We denote the probability of belonging to component  $k$  by  $w_k$ , such that  $\sum_{k=1}^K w_k = 1$ .  $\boldsymbol{\mu}_k$  is the mean of the  $k$ -th component and it is assumed that each component has the same covariance matrix  $\boldsymbol{\Sigma}_b$ . This constraint is needed to avoid infinite likelihoods (Böhning, 1999). We write  $\mathbf{w} = (w_1, \dots, w_{K-1})^T$ , the vector of  $K - 1$  component probabilities, and  $\boldsymbol{\mu} = (\boldsymbol{\mu}_1^T, \dots, \boldsymbol{\mu}_K^T)^T$ , the vector of all component means. We introduce  $\mathbf{b}_i$  and  $\alpha_i = k$ , ( $k = 1, \dots, K$ ), as the  $i$ -th subject's random effects following the mixture distribution and the  $k$ -th component of the mixture from which  $\mathbf{b}_i$  is sampled, respectively. The distribution of  $\alpha_i$  is then described by  $P(\alpha_i = k) = w_k$  and, given  $\alpha_i = k$ ,  $\mathbf{b}_{ik} \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_b)$ . Thus,  $\mathbf{b}_i = \sum_{k=1}^K I(\alpha_i = k) \mathbf{b}_{ik}$ , where  $I(\alpha_i = k)$  is the indicator of belonging to component  $k$ . For  $n$  subjects,  $\mathbf{b} = (\mathbf{b}_1^T, \dots, \mathbf{b}_n^T)^T$  and  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)^T$ .

Now we estimate and make inferences on the parameters  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \boldsymbol{\phi}^T, \text{Vec}(\boldsymbol{\Sigma}_b)^T, \boldsymbol{\mu}^T, \mathbf{w}^T, \boldsymbol{\psi}^T, \boldsymbol{\gamma}^T)^T$  and the baseline cumulative hazard functions with  $S$  strata,  $\boldsymbol{\Lambda}(t) = (\Lambda_1(t), \dots, \Lambda_S(t))^T$ , where  $\Lambda_s(t) = \int_0^t \lambda_s(u) du$ ,  $s = 1, \dots, S$ . The parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\phi}$  are from the longitudinal model,  $\boldsymbol{\psi}$  and  $\boldsymbol{\gamma}$  are from the hazard model, and,  $\boldsymbol{\mu}$ ,  $\mathbf{w}$ , and  $\boldsymbol{\Sigma}_b$  are associated with the random effects.  $\text{Vec}(\cdot)$  operator creates a column vector from a matrix by stacking the diagonal and upper-triangle elements of the matrix. The likelihood function (4.3) of the complete data  $(\mathbf{Y}, \mathbf{V}, \mathbf{b}, \boldsymbol{\alpha})$  and the full likelihood function (4.4) of the observed data  $(\mathbf{Y}, \mathbf{V})$  for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  have the following forms respectively,

$$\begin{aligned}
& L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}, \boldsymbol{\alpha}) \\
&= \prod_{i=1}^n \prod_{k=1}^K \left[ \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_{ik}) - B(\boldsymbol{\beta}; \mathbf{b}_{ik})}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \\
&\quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k) \right\} \times w_k \left. \right]^{I(\alpha_i=k)}
\end{aligned}$$

$$\text{and} \quad L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = \sum_{\boldsymbol{\alpha}} \int_{\mathbf{b}} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}, \boldsymbol{\alpha}) d\mathbf{b}.$$

The proposed estimation method is to calculate the maximum likelihood estimates for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda}(t))$  over a set of  $\boldsymbol{\theta}$  and  $\boldsymbol{\Lambda}(t)$ . We let each  $\Lambda_s(t)$  of  $\boldsymbol{\Lambda}(t)$ ,  $s = 1, \dots, S$ , be a non-decreasing and right-continuous step function with jumps only at the observed failure times belonging to stratum  $s$ .

EM-algorithm is used for calculating the maximum likelihood estimates. In the EM-algorithm,  $\mathbf{b}_i$  and  $\alpha_i$  are considered as missing data for  $i = 1, \dots, n$ . Therefore, the M-step solves the conditional score equations from complete data given observations, where the conditional expectation can be evaluated in E-step. The procedure involves iterating between the following two steps until convergence is achieved: at the  $m$ -th iteration,

(1) E-step Calculate the conditional expectations of some known functions of  $\mathbf{b}_i$  and  $\alpha_i$ , needed in the next M-step, for subject  $i$  with  $S_i = s$  given observations and the current estimate  $(\boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)})$ . The conditional expectation is calculated using the Gauss-Hermite Quadrature numerical approximation, denoted as  $E[q(\mathbf{b}_i, \alpha_i) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}]$  for a

known function  $q(\mathbf{b}_i, \alpha_i)$ .

(2) M-step After differentiating the conditional expectation of complete data log-likelihood function given observations and the current estimate  $(\boldsymbol{\theta}^{(m)}, \boldsymbol{\Lambda}^{(m)})$ , the updated estimator  $(\boldsymbol{\theta}^{(m+1)}, \boldsymbol{\Lambda}^{(m+1)})$  can be obtained as follows:  $(\boldsymbol{\beta}^{(m+1)}, \boldsymbol{\phi}^{(m+1)})$  solves the conditional expectation of complete data log-likelihood score equation using one-step Newton-Raphson iteration; For the covariance matrix of random effects,

$$\boldsymbol{\Sigma}_b^{(m+1)} = \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \sum_{k=1}^K \mathbb{E} [I(\alpha_i = k)(\mathbf{b}_{ik} - \boldsymbol{\mu}_k)(\mathbf{b}_{ik} - \boldsymbol{\mu}_k)^T | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] I(S_i = s);$$

For the  $k$ -th mixture component ( $k = 1, \dots, K$ ),

$$\boldsymbol{\mu}_k^{(m+1)} = \frac{\sum_{i=1}^n \sum_{s=1}^S \mathbb{E} [I(\alpha_i = k) \mathbf{b}_{ik} | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] I(S_i = s)}{\sum_{i=1}^n \sum_{s=1}^S \mathbb{E} [I(\alpha_i = k) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] I(S_i = s)}$$

and  $w_k^{(m+1)} = \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \mathbb{E} [I(\alpha_i = k) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] I(S_i = s);$

$(\boldsymbol{\psi}^{(m+1)}, \boldsymbol{\gamma}^{(m+1)})$  solves the partial likelihood score equation from the full data using one-step Newton-Raphson iteration,

$$\begin{aligned} & \sum_{i=1}^n \sum_{s=1}^S \Delta_i \left\{ \begin{pmatrix} E[(\tilde{\mathbf{Z}}_i^T(V_i) \circ \mathbf{b}_i) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] \\ \mathbf{Z}_i \end{pmatrix} \right. \\ & \quad \left. \frac{\sum_{l: V_l \geq V_i} \left( \begin{pmatrix} E[(\tilde{\mathbf{Z}}_l^T(V_i) \circ \mathbf{b}_l) \exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma}\} | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] \\ E[\mathbf{Z}_l(V_i) \exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma}\} | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] \end{pmatrix} I(S_l = s) \right)}{E[\exp\{\tilde{\mathbf{Z}}_l(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i)\boldsymbol{\gamma}\} | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] I(S_l = s)} \right\} I(S_i = s) \\ & = \mathbf{0}; \end{aligned}$$

$\Lambda_s^{(m+1)}$  is obtained as an empirical function with jumps only at the observed failure time,

$$\Lambda_s^{(m+1)}(t) = \sum_{i: V_i \leq t} \frac{\Delta_i I(S_i = s)}{\sum_{l: V_l \geq V_i} E \left[ \exp \left\{ \tilde{\mathbf{Z}}_l(V_i) (\boldsymbol{\psi}^{(m+1)} \circ \mathbf{b}_l) + \mathbf{Z}_l(V_i) \boldsymbol{\gamma}^{(m+1)} \right\} | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)} \right] I(S_l = s)}.$$

The expressions of the conditional expectation and the conditional score equations calculated in the E- and M-steps for continuous longitudinal outcomes following a normal distribution and binary longitudinal outcomes with survival time are given respectively in Appendices A.1 and A.2.

The observed information matrix via Louis (1982) formula is adopted to obtain the variance estimate for  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}(t))$ . The variance of  $\sqrt{n} \widehat{\boldsymbol{\theta}}$  is asymptotically equal to the corresponding sub-matrix of the inverse of the calculated observed information matrix.

### 4.2.3 EM algorithm – examples

#### 4.2.3.1 Continuous longitudinal data with Normal distribution and survival time

(1) E-step : For continuous longitudinal outcomes following a normal distribution and survival time, we calculate the conditional expectation of  $q(\mathbf{b}_i, \alpha_i)$  for subject  $i$  with  $S_i = s$  given the observations and the current estimate  $(\boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)})$  for some known function  $q(\cdot)$ . The conditional expectation denoted by  $E[q(\mathbf{b}_i, \alpha_i) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}]$  can be expressed as the following: Given the current estimate  $(\boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)})$ ,

$$E[q(\mathbf{b}_i, \alpha_i) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}] = \frac{\sum_{\alpha_i=1}^K C_{\alpha_i} \int_{\mathbf{z}_{\alpha_i}} q(R(\mathbf{z}_{\alpha_i})) \kappa(\mathbf{z}_{\alpha_i}) \exp\{-\mathbf{z}_{\alpha_i}^T \mathbf{z}_{\alpha_i}\} d\mathbf{z}_{\alpha_i}}{\sum_{a=1}^K C_a \int_{\mathbf{z}_a} \kappa(\mathbf{z}_a) \exp\{-\mathbf{z}_a^T \mathbf{z}_a\} d\mathbf{z}_a}, \quad (4.5)$$

where

$$R(\mathbf{z}_{\alpha_i}) = \left( \frac{1}{\sigma_y^2} \sum_{j=1}^{n_i} \tilde{\mathbf{X}}_{ij}^T \tilde{\mathbf{X}}_{ij} + (\boldsymbol{\Sigma}_b^{(m)})^{-1} \right)^{-1} \left[ \sum_{j=1}^{n_i} \frac{1}{\sigma_y^2} (y_{ij} - \mathbf{X}_{ij} \boldsymbol{\beta}) \tilde{\mathbf{X}}_{ij}^T + \Delta_i (\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) \right]$$

$$\begin{aligned}
& + (\Sigma_b^{(m)})^{-1} \boldsymbol{\mu}_{\alpha_i} \Big] + \sqrt{2} \left[ \frac{1}{\sigma_y^2} \sum_{j=1}^{n_i} \tilde{\mathbf{X}}_{ij}^T \tilde{\mathbf{X}}_{ij} + (\Sigma_b^{(m)})^{-1} \right]^{-\frac{1}{2}} \mathbf{z}_{\alpha_i}, \\
\kappa(\mathbf{z}_{\alpha_i}) &= \exp \left\{ - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{\tilde{\mathbf{Z}}_i(u) (\boldsymbol{\psi}^{(m)} \circ R(\mathbf{z}_{\alpha_i})) + \mathbf{Z}_i(u) \boldsymbol{\gamma}^{(m)}} d\Lambda_s^{(m)}(u) \right\}, \\
C_{\alpha_i} &= \exp \left\{ \frac{1}{2} \left[ \sum_{j=1}^{n_i} \frac{1}{\sigma_y^2} (y_{ij} - \mathbf{X}_{ij} \boldsymbol{\beta}) \tilde{\mathbf{X}}_{ij}^T + \Delta_i (\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) + (\Sigma_b^{(m)})^{-1} \boldsymbol{\mu}_{\alpha_i} \right]^T \right. \\
& \quad \times \left( \frac{1}{\sigma_y^2} \sum_{j=1}^{n_i} \tilde{\mathbf{X}}_{ij}^T \tilde{\mathbf{X}}_{ij} + (\Sigma_b^{(m)})^{-1} \right)^{-1} \times \left[ \sum_{j=1}^{n_i} \frac{1}{\sigma_y^2} (y_{ij} - \mathbf{X}_{ij} \boldsymbol{\beta}) \tilde{\mathbf{X}}_{ij}^T \right. \\
& \quad \left. \left. + \Delta_i (\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) + (\Sigma_b^{(m)})^{-1} \boldsymbol{\mu}_{\alpha_i} \right] \right. \\
& \quad \left. - \frac{1}{2} \boldsymbol{\mu}_{\alpha_i}^T (\Sigma_b^{(m)})^{-1} \boldsymbol{\mu}_{\alpha_i} + \log w_{\alpha_i} \right\}
\end{aligned}$$

is a constant,

$(\Sigma_b^{(m)})^{\frac{1}{2}}$  is an unique non-negative square root of  $\Sigma_b^{(m)}$  (i.e.  $(\Sigma_b^{(m)})^{\frac{1}{2}} \times (\Sigma_b^{(m)})^{\frac{1}{2}} = \Sigma_b^{(m)}$ ),

and  $\mathbf{z}_{\alpha_i}$  follows a multivariate Gaussian distribution with mean zero.

(2) M-step : Since normal distribution has a dispersion parameter  $\phi$  as  $\sigma_y^2$ , we estimate  $\boldsymbol{\beta}^{(m+1)}$  and  $\sigma_y^2$  in longitudinal process.  $\boldsymbol{\beta}^{(m+1)}$  is the linear regression coefficients of regressing  $\{\mathbf{Y}_i - \mathbb{E}[\tilde{\mathbf{X}}_i \mathbf{b}_i | \boldsymbol{\theta}^{(m)}, \Lambda^{(m)}], i = 1, \dots, n\}$  on  $\{\mathbf{X}_i, i = 1, \dots, n\}$ , where  $\mathbf{X}_i = (\mathbf{X}_{i1}^T, \dots, \mathbf{X}_{in_i}^T)^T$  and  $\tilde{\mathbf{X}}_i = (\tilde{\mathbf{X}}_{i1}^T, \dots, \tilde{\mathbf{X}}_{in_i}^T)^T$ .

$$(\sigma_y^2)^{(m+1)} = \frac{\sum_{i=1}^n [D_i^T D_i + \mathbb{E}[(\tilde{\mathbf{X}}_i \mathbf{b}_i)^2 | \boldsymbol{\theta}^{(m)}, \Lambda^{(m)}] - (\mathbb{E}[\tilde{\mathbf{X}}_i \mathbf{b}_i | \boldsymbol{\theta}^{(m)}, \Lambda^{(m)}])^2]}{\sum_{i=1}^n n_i},$$

where  $D_i = \mathbf{Y}_i - \mathbf{X}_i \boldsymbol{\beta}^{(m+1)} - \mathbb{E}[\tilde{\mathbf{X}}_i \mathbf{b}_i | \boldsymbol{\theta}^{(m)}, \Lambda^{(m)}]$ .  $\Sigma_b^{(m+1)}$ ,  $\boldsymbol{\mu}^{(m+1)}$ ,  $\mathbf{w}^{(m+1)}$ ,  $(\boldsymbol{\psi}^{(m+1)}, \boldsymbol{\gamma}^{(m+1)})$ , and  $\Lambda_s^{(m+1)}$  have the same expressions as in Section 4.2.2.

#### 4.2.3.2 Binary longitudinal data and survival time

(1) E-step : For binary longitudinal outcomes and survival time, given the current estimate  $(\boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)})$ , the conditional expectation denoted by  $E[q(\mathbf{b}_i, \alpha_i) | \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)}]$



can be expressed as in (4.5), where

$$\begin{aligned}
R(\mathbf{z}_{\alpha_i}) &= \Sigma_b^{(m)} \left[ \sum_{j=1}^{n_i} y_{ij} \tilde{\mathbf{X}}_{ij}^T + \Delta_i(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) \right] + \boldsymbol{\mu}_{\alpha_i} + \sqrt{2}(\Sigma_b^{(m)})^{\frac{1}{2}} \mathbf{z}_{\alpha_i}, \\
\kappa(\mathbf{z}_{\alpha_i}) &= \exp \left\{ - \sum_{j=1}^{n_i} \log \left( 1 + e^{\mathbf{X}_{ij} \boldsymbol{\beta}^{(m)} + \tilde{\mathbf{X}}_{ij} R(\mathbf{z}_{\alpha_i})} \right) \right. \\
&\quad \left. - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi}^{(m)} \circ R(\mathbf{z}_{\alpha_i})) + \mathbf{Z}_i(u) \boldsymbol{\gamma}^{(m)}} d\Lambda_s^{(m)}(u) \right\}, \text{ and} \\
C_{\alpha_i} &= \exp \left\{ \frac{1}{2} \left[ \Sigma_b^{(m)} \left( \sum_{j=1}^{n_i} y_{ij} \tilde{\mathbf{X}}_{ij}^T + \Delta_i(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) \right) + \boldsymbol{\mu}_{\alpha_i} \right]^T \times (\Sigma_b^{(m)})^{-1} \right. \\
&\quad \times \left[ \Sigma_b^{(m)} \left( \sum_{j=1}^{n_i} y_{ij} \tilde{\mathbf{X}}_{ij}^T + \Delta_i(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi}^{(m)}) \right) + \boldsymbol{\mu}_{\alpha_i} \right] \\
&\quad \left. - \frac{1}{2} \boldsymbol{\mu}_{\alpha_i}^T (\Sigma_b^{(m)})^{-1} \boldsymbol{\mu}_{\alpha_i} + \log w_{\alpha_i} \right\}
\end{aligned}$$

is a constant.

(2) M-step : Since the parameter  $\boldsymbol{\phi}$  is set to 1 for logistic distribution, we estimate only  $\boldsymbol{\beta}$  in the longitudinal process.  $\boldsymbol{\beta}^{(m+1)}$  solves the conditional expectation of complete data log-likelihood score equation, using one-step Newton-Raphson iteration,

$$\sum_{i=1}^n \sum_{j=1}^{n_i} \left( y_{ij} - \sum_{s=1}^S \mathbb{E} \left[ \frac{\exp\{\mathbf{X}_{ij} \boldsymbol{\beta}^{(m+1)} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i\}}{1 + \exp\{\mathbf{X}_{ij} \boldsymbol{\beta}^{(m+1)} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i\}} \middle| \boldsymbol{\theta}^{(m)}, \Lambda_s^{(m)} \right] I(S_i = s) \right) \mathbf{X}_{ij}^T = \mathbf{0}.$$

$\Sigma_b^{(m+1)}, \boldsymbol{\mu}^{(m+1)}, \mathbf{w}^{(m+1)}, \boldsymbol{\psi}^{(m+1)}, \boldsymbol{\gamma}^{(m+1)}$ , and  $\Lambda_s^{(m+1)}$  have the same expressions as in Section 4.2.2.

### 4.3 Asymptotic Properties

To study the asymptotic properties of the proposed estimator  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}}(t))$  with  $\hat{\boldsymbol{\theta}} = (\hat{\boldsymbol{\beta}}^T, \hat{\boldsymbol{\phi}}^T, \text{Vec}(\hat{\Sigma}_b)^T, \boldsymbol{\mu}^T, \mathbf{w}^T, \hat{\boldsymbol{\psi}}^T, \hat{\boldsymbol{\gamma}}^T)^T$  and  $\hat{\boldsymbol{\Lambda}}(t) = (\hat{\Lambda}_1(t), \dots, \hat{\Lambda}_S(t))^T$ , we assume the following conditions.

(A1) The true parameter  $\boldsymbol{\theta}_0 = (\boldsymbol{\beta}_0^T, \boldsymbol{\phi}_0^T, \text{Vec}(\Sigma_{b0})^T, \boldsymbol{\mu}^T, \mathbf{w}^T, \boldsymbol{\psi}_0^T, \boldsymbol{\gamma}_0^T)^T$  belongs to a known

compact set  $\Theta$  which lies in the interior of the domain for  $\theta$ .

- (A2) The distribution of random effects  $\mathbf{b}_i^*$  is a mixture of a finite number of  $d_b$ -dimensional multivariate normal distributions with means  $\boldsymbol{\mu} = (\boldsymbol{\mu}_1^T, \dots, \boldsymbol{\mu}_K^T)^T$  and a common covariance matrix  $\boldsymbol{\Sigma}_b$ . i.e.  $\mathbf{b}_i^* \sim \sum_{k=1}^K w_k \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_b)$ , where  $K$  is the number of mixture components.
- (A3) The true baseline hazard rate function  $\boldsymbol{\lambda}_0(t) = (\lambda_{10}(t), \dots, \lambda_{S0}(t))$  is continuous and positive in  $[0, \tau]$ , where  $\tau$  is the time of study end.
- (A4) For the censoring time  $C$ ,  $P(C \geq \tau | \mathbf{Z}, \tilde{\mathbf{Z}}, \mathbf{X}, \tilde{\mathbf{X}}) = P(C = \tau | \mathbf{Z}, \tilde{\mathbf{Z}}, \mathbf{X}, \tilde{\mathbf{X}}) > 0$ .
- (A5) For the number of observed longitudinal measurements per subject  $n_N$ ,  $P(n_N > d_b | \mathbf{X}, \tilde{\mathbf{X}}) > 0$  with probability one, and  $P(n_N \leq n_0) = 1$  for some integer  $n_0$ .
- (A6) Both  $\mathbf{X}^T \mathbf{X}$  and  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  are full rank with positive probability. Moreover, if there exist constant vectors  $\mathbf{c}_1$  and  $\mathbf{c}_2$  such that, with positive probability, for any  $t$ ,  $\mathbf{Z}(t)\mathbf{c}_1 = \alpha_0(t)$  and  $\tilde{\mathbf{Z}}(t) \circ \mathbf{c}_2 = 0$  for a deterministic function  $\alpha_0(t)$ , then  $\mathbf{c}_1 = 0$ ,  $\mathbf{c}_2 = 0$ , and  $\alpha_0(t) = 0$ .

Assumption (A4) means that, by the end of the study, some proportion of the subjects will still be alive and censored at the study end time  $\tau$ , and thus the maximum right censoring time is equal to  $\tau$ . Assumption (A5) implies that some proportion of the subjects have at least  $d_b$  longitudinal observations, and there exists an integer  $n_0$  such that all subjects have a finite number of longitudinal observations which are not larger than  $n_0$ . Consistency and asymptotic distribution of the proposed estimator are summarized in the following two theorems. We will present outlines of the proofs here. The detailed technical proofs are given in Section 4.4.

**Theorem 4.1.** *Under the assumptions (A1)~(A6), as  $n \rightarrow \infty$ , the maximum likelihood estimator  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\Lambda}}(t))$  is consistent under the product norm of the Euclidean distance and*

the supreme norm on  $[0, \tau]$ . That is,  $\|\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} \|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$ , a.s., where  $\|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| = \sum_{s=1}^S |\widehat{\Lambda}_s(t) - \Lambda_{s0}(t)|$ .

Consistency in Theorem 4.1 can be proved by verifying the following three steps: First, we show that the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  exists. Second, we show that, with probability one,  $\widehat{\Lambda}_s(\tau)$ ,  $s = 1, \dots, S$ , are bounded as  $n \rightarrow \infty$ . Third, given that the second step is true, by Helly's selection theorem (van der Vaart, 1998), we can choose a subsequence of  $\widehat{\Lambda}_s(t)$  such that  $\widehat{\Lambda}_s(t)$  weakly converges to some right-continuous monotone function  $\Lambda_s^*(t)$  with probability one. Also, for any sub-sequence, we can find a further sub-sequence, still denoted as  $\widehat{\boldsymbol{\theta}}$ , such that  $\widehat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$ . Using empirical process formulation and relevant Donsker properties with parameter identifiability, we can show that  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$ ,  $s = 1, \dots, S$ . Once the three steps are completed, we can conclude that, with probability one,  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\widehat{\boldsymbol{\Lambda}}_s(t)$  converges to  $\Lambda_{s0}(t)$  in  $[0, \tau]$ ,  $s = 1, \dots, S$ . Moreover, since  $\Lambda_{s0}(t)$  is right-continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$  almost surely.

**Theorem 4.2.** *Under the assumptions (A1)~(A6), as  $n \rightarrow \infty$ ,  $\sqrt{n}((\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T, (\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t))^T)^T$  weakly converges to a Gaussian random element in  $R^{d_\theta} \times \ell^\infty[0, \tau] \times \dots \times \ell^\infty[0, \tau]$ , and the estimator  $\widehat{\boldsymbol{\theta}}$  is asymptotically efficient, where  $d_\theta$  is the dimension of  $\boldsymbol{\theta}$  and  $\ell^\infty[0, \tau]$  is the normed space containing all the bounded functions in  $[0, \tau]$ .*

Once consistency holds, the conditions of Theorem 3.3.1 in van der Vaart and Wellner (1996), which implies the asymptotic normality in Theorem 4.2, are verified via the tools of empirical processes. These conditions are restated in Theorem 4 of Parner (1998). The smooth conditions in Theorem 4 of Parner (1998) can be verified using the regularity of the log-likelihood function in terms of model parameters and the Donsker properties of the score operators. In particular, in the invertibility condition of the

information operator in Theorem 4 of Parner (1998), the verification of the one-to-one property of the information operator is specific to our proposed models and requires non-trivial work. Therefore, by Theorem 3.3.1 of van der Vaart and Wellner (1996),  $\sqrt{n}((\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T, (\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t))^T)^T$  weakly converges to a Gaussian process, and by Proposition 3.3.1 in Bickel *et al.* (1993),  $\widehat{\boldsymbol{\theta}}$  is an efficient estimator for  $\boldsymbol{\theta}_0$ .

#### 4.4 Technical Details – Proofs for Asymptotic Properties

In this section, we present the detailed technical proofs for the asymptotic properties of the proposed estimator  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}(t))$  with  $\widehat{\boldsymbol{\theta}} = (\widehat{\boldsymbol{\beta}}^T, \widehat{\boldsymbol{\phi}}^T, \text{Vec}(\widehat{\boldsymbol{\Sigma}}_b)^T, \widehat{\boldsymbol{\mu}}^T, \widehat{\boldsymbol{w}}^T, \widehat{\boldsymbol{\psi}}^T, \widehat{\boldsymbol{\gamma}}^T)^T$  and  $\widehat{\boldsymbol{\Lambda}}(t) = (\widehat{\Lambda}_1(t), \dots, \widehat{\Lambda}_S(t))^T$ . Meanwhile, the supplementary proofs needed to prove the asymptotic properties are provided in Section 4.4.3. We use  $\mathbf{b}$  to denote random effects instead of  $\mathbf{b}^*$  for convenience in all proofs. From the following full likelihood function of observed data  $(\mathbf{Y}, \mathbf{V})$  for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda})$ ,

$$\begin{aligned}
& L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) \\
&= \sum_{\boldsymbol{\alpha}} \int_{\mathbf{b}} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}, \boldsymbol{\alpha}) d\mathbf{b} \\
&= \prod_{i=1}^n \left( \sum_{\alpha_i=1}^K \int_{\mathbf{b}_1} \cdots \int_{\mathbf{b}_K} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}_i, V_i, \mathbf{b}_{\alpha_i}, \alpha_i) d\mathbf{b}_1 \cdots d\mathbf{b}_K \right) \\
&= \prod_{i=1}^n \left( \sum_{\alpha_i=1}^K \int_{\mathbf{b}_{\alpha_i}} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}_i, V_i, \mathbf{b}_{\alpha_i}, \alpha_i) d\mathbf{b}_{\alpha_i} \right) \\
&= \prod_{i=1}^n \left( \sum_{\alpha_i=1}^K \prod_{k=1}^K \left[ \int_{\mathbf{b}_k} \left( \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_{ik}) - B(\boldsymbol{\beta}; \mathbf{b}_{ik})}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \right. \right. \\
&\quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\widetilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \widetilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\quad \left. \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k) \right\} \times w_k \right) d\mathbf{b}_k \Bigg]^{I(\alpha_i=k)} \right),
\end{aligned}$$

we have the observed log-likelihood function

$$\begin{aligned}
& \sum_{i=1}^n \log \left( \sum_{\alpha_i=1}^K \prod_{k=1}^K \left[ \int_{\mathbf{b}_k} \left( \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_{ik}) - B(\boldsymbol{\beta}; \mathbf{b}_{ik})}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \right. \\
& \quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
& \quad \quad \left. \left. \left. - \int_0^{V_i} \exp \left\{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \right\} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
& \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k) \right\} \times w_k \left. \right]^{I(\alpha_i=k)} \Bigg).
\end{aligned}$$

Then, we obtain the following modified object function by replacing  $\lambda_s(V_i)$  with  $\Lambda_s\{V_i\}$  in the above expression where  $\Lambda_s\{V_i\}$  is the jump size of  $\Lambda_s(t)$  at the observed time  $V_i$  with  $\Delta_i = 1$ ,

$$\begin{aligned}
& l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \\
& = \sum_{i=1}^n \log \left( \sum_{\alpha_i=1}^K \prod_{k=1}^K \left[ \int_{\mathbf{b}_k} \left( \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_{ik}) - B(\boldsymbol{\beta}; \mathbf{b}_{ik})}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right. \right. \\
& \quad \times \left( \prod_{s=1}^S \left[ \Lambda_s\{V_i\}^{\Delta_i} \exp \left\{ \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
& \quad \quad \left. \left. \left. - \int_0^{V_i} \exp \left\{ \tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_{ik}) + \mathbf{Z}_i(u)\boldsymbol{\gamma} \right\} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
& \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_{ik} - \boldsymbol{\mu}_k) \right\} \times w_k \left. \right]^{I(\alpha_i=k)} \Bigg), \tag{4.6}
\end{aligned}$$

and  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  maximizes  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  over the space  $\{(\boldsymbol{\theta}, \boldsymbol{\Lambda}) : \boldsymbol{\theta} \in \Theta, \boldsymbol{\Lambda} \in \mathbb{W}_n \times \mathbb{W}_n \cdots \times \mathbb{W}_n\}$ , where  $\mathbb{W}_n$  consists of all the right-continuous step functions only; that is,  $\boldsymbol{\Lambda} = (\Lambda_1, \dots, \Lambda_S)^T, s = 1, \dots, S, \Lambda_s \in \mathbb{W}_n$ . For the proofs of both Theorem 4.1 and Theorem 4.2, the modified object function is used in the place of the observed log-likelihood function.

#### 4.4.1 Proof of consistency

Consistency can be proved by verifying the following three steps: First, we show the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  exists. Second, we show that, with probability one,  $\widehat{\Lambda}_s(\tau)$ ,  $s = 1, \dots, S$ , are bounded as  $n \rightarrow \infty$ . Third, if the second step is true, by Helly's selection theorem (p9 of van der Vaart (1998)), we can choose a subsequence of  $\widehat{\Lambda}_s$  such that  $\widehat{\Lambda}_s$  weakly converges to some right-continuous monotone function  $\Lambda_s^*$  with probability one; that is, the measure given by  $\mu_s([0, t]) = \widehat{\Lambda}_s(t)$  for  $t \in [0, \tau]$  weakly converges to the measure given by  $\mu_s^*([0, t]) = \Lambda_s^*(t)$ . By choosing a sub-sequence, we can further assume  $\widehat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$ . Thus, in this third step, we show  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$ ,  $s = 1, \dots, S$ .

Once the three steps are completed, we can conclude that, with probability one,  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\widehat{\Lambda}_s$  converges to  $\Lambda_{s0}$  in  $[0, \tau]$ ,  $s = 1, \dots, S$ . However, since  $\Lambda_{s0}$  is continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\widehat{\Lambda}(t) - \Lambda_0(t)\| \rightarrow 0$  almost surely. Then, the proof of Theorem 4.1 will be done.

In the first step, we will show the existence of the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$ . Since  $\boldsymbol{\theta}$  belongs to a compact set  $\Theta$  by the assumption (A1), it is sufficient to show that  $\Lambda_s\{V_i\}$ , the jump size of  $\Lambda_s$  at  $V_i$  for which  $\Delta_i = 1$ , is finite. Since, for each subject  $i$  with  $\Delta_i = 1$ ,

$$\begin{aligned} \Lambda_s\{V_i\} \exp \left\{ - \int_0^{V_i} \exp \{ \widetilde{\mathbf{Z}}_i(t)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}_i(t)\boldsymbol{\gamma} \} d\Lambda_s(t) \right\} \\ \leq \exp \left\{ - 2(\widetilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}) \right\} (\Lambda_s\{V_i\})^{-1}, \end{aligned}$$

we have that, from (4.6),

$$\begin{aligned} l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \\ \leq \sum_{i=1}^n \log \left( \sum_{\alpha_i=1}^K \prod_{k=1}^K \left[ \int_{\mathbf{b}_k} \left( \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_k) - B(\boldsymbol{\beta}; \mathbf{b}_k)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \right) \right] \right) \end{aligned}$$

$$\begin{aligned} & \times \left( \prod_{s=1}^S \left[ \left( \Lambda_s \{V_i\} \right)^{-\Delta_i} \exp \left\{ -\Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_k) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right\} \right]^{I(S_i=s)} \right) \\ & \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_k - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_k - \boldsymbol{\mu}_k) \right\} \times w_k \Big) d\mathbf{b}_k \Big]^{I(\alpha_i=k)} \Big). \end{aligned}$$

Thus, if  $\Lambda_s \{V_i\} \rightarrow \infty$  for some  $i$  with  $\Delta_i = 1$ , then  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \rightarrow -\infty$ , which is contradictory to that  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  is bounded. Therefore, we conclude that  $\Lambda_s \{\cdot\}$ , the jump size of  $\Lambda_s$  for stratum  $s$ , must be finite. By the conclusion and the assumption (A1), the maximum likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  exists.

In the second step, we will show that  $\widehat{\Lambda}_s(\tau)$  is bounded as  $n$  goes to infinity with probability one. We define  $\widehat{\zeta}_s = \log \widehat{\Lambda}_s(\tau)$  and rescale  $\widehat{\Lambda}_s$  by the factor  $e^{\widehat{\zeta}_s}$ . Then, we let  $\widetilde{\Lambda}_s$  denote the rescaled function; that is,  $\widetilde{\Lambda}_s(t) = \widehat{\Lambda}_s(t)/\widehat{\Lambda}_s(\tau) = \widehat{\Lambda}_s(t)e^{-\widehat{\zeta}_s}$ . thus,  $\widetilde{\Lambda}_s(\tau) = 1$ . To prove this second step, it is sufficient to show  $\widehat{\zeta}_s$  is bounded. After some algebra in (4.6), we obtain that, for any  $\boldsymbol{\Lambda} \in \mathbb{W} \times \mathbb{W} \cdots \times \mathbb{W}$ ,

$$\begin{aligned} & n^{-1} l_n(\widehat{\boldsymbol{\theta}}, \boldsymbol{\Lambda}) \\ &= \frac{1}{n} \sum_{i=1}^n \left[ \sum_{j=1}^{n_i} \left( \frac{Y_{ij} \mathbf{X}_{ij} \widehat{\boldsymbol{\beta}}}{A(D_i(t_j; \widehat{\phi}))} + C(Y_{ij}; D_i(t_j; \widehat{\phi})) \right) + \sum_{s=1}^S \Delta_i I(S_i=s) \log \Lambda_s \{V_i\} + \Delta_i \mathbf{Z}_i(V_i) \boldsymbol{\gamma} \right. \\ & \quad - \frac{1}{2} \log \{ (2\pi)^{d_b} |\widehat{\boldsymbol{\Sigma}}_b| \} - \frac{1}{2} \log |\widehat{\boldsymbol{\Sigma}}_b| + \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\ & \quad \times \int_{\mathbf{b}_{\alpha 0}} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right. \right. \\ & \quad \left. \left. \left. - \sum_{s=1}^S I(S_i=s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\Lambda_s(t) \right\} \right] d\mathbf{b}_{\alpha 0} \right] \Big], \end{aligned}$$

where

$$\begin{aligned} \mathbf{M}_{i\alpha} &= \left[ \left( \sum_{j=1}^{n_i} \frac{Y_{ij} \widetilde{\mathbf{X}}_{ij}}{A(D_i(t_j; \widehat{\phi}))} + \Delta_i (\widetilde{\mathbf{Z}}_i(V_i) \circ \widehat{\boldsymbol{\psi}}^T) \right) \widehat{\boldsymbol{\Sigma}}_b^{1/2} + \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1/2} \right]^T, \\ \mathbf{b}_{\alpha 0} &= \boldsymbol{\Sigma}_b^{-1/2} \mathbf{b}_\alpha - \mathbf{M}_{i\alpha}, \end{aligned}$$

and

$$Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}}) = (\widetilde{\mathbf{Z}}_i(t) \circ \widehat{\boldsymbol{\psi}}^T) \widehat{\boldsymbol{\Sigma}}_b^{1/2} \mathbf{b}_{\alpha 0} + \mathbf{Z}_i(t) \boldsymbol{\gamma} + (\widetilde{\mathbf{Z}}_i(t) \circ \widehat{\boldsymbol{\psi}}^T) \widehat{\boldsymbol{\Sigma}}_b^{1/2} \mathbf{M}_{i\alpha}.$$

Thus, since  $0 \leq n^{-1}l_n(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}) - n^{-1}l_n(\widehat{\boldsymbol{\theta}}, \widetilde{\boldsymbol{\Lambda}})$  where  $\widehat{\boldsymbol{\Lambda}} = e^{\widehat{\boldsymbol{\xi}}} \circ \widetilde{\boldsymbol{\Lambda}}$ , it follows that

$$\begin{aligned}
0 \leq & \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S I(S_i = s) \Delta_i \left( \log e^{\widehat{\zeta}_s} \widetilde{\Lambda}_s \{V_i\} - \log \widetilde{\Lambda}_s \{V_i\} \right) \\
& + \frac{1}{n} \sum_{i=1}^n \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right. \\
& \quad \quad \left. \left. - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \right] \Bigg] \\
& - \frac{1}{n} \sum_{i=1}^n \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right. \\
& \quad \quad \left. \left. - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \right] \Bigg]. \quad (4.7)
\end{aligned}$$

According to the assumption (A3), there exist some positive constants  $C1$ ,  $C2$  and  $C3$  such that  $|Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})| \leq C_1 \|\mathbf{b}_{\alpha 0}\| + C_2 \|\mathbf{Y}_i\| + C_3$ . By denoting  $\mathbf{b}_{\alpha 0}$  as a vector of variables following a standard multivariate normal distribution, from concavity of the logarithm function, in the third term of (4.7),

$$\begin{aligned}
& \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \Bigg] \Bigg] \\
& = \frac{d_b}{2} \log(2\pi) + \log \left[ E_\alpha \left[ \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \times E_{\mathbf{b}_{\alpha 0}|\alpha} \left[ \exp \left\{ -\sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \right] \Bigg] \Bigg] \\
& = \frac{d_b}{2} \log(2\pi) + \log \left[ E_{\alpha, \mathbf{b}_0} \left[ \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \times \exp \left\{ -\sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \Bigg] \Bigg]
\end{aligned}$$



$$\begin{aligned}
&\geq \frac{d_b}{2} \log(2\pi) + \log \left[ E_{\alpha, \mathbf{b}_0} \left[ \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} \right. \right. \right. \\
&\quad \left. \left. \left. - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - e^{C_1 \|\mathbf{b}_{\alpha 0}\| + C_2 \|\mathbf{Y}_i\| + C_3} \right\} \right] \right] \\
&\geq \frac{d_b}{2} \log(2\pi) + \log \left[ E_{\alpha, \mathbf{b}_0} \left[ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - e^{C_1 \|\mathbf{b}_{\alpha 0}\| + C_2 \|\mathbf{Y}_i\| + C_3} \right] \right] \\
&= -e^{C_2 \|\mathbf{Y}_i\| + C_4} - C_5,
\end{aligned}$$

where  $C_4$  and  $C_5$  are positive constants. Then, since it is easily verified that  $E_{\alpha, \mathbf{b}_0} \left[ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - e^{C_1 \|\mathbf{b}_{\alpha 0}\| + C_2 \|\mathbf{Y}_i\| + C_3} \right] < \infty$ , by the strong law of large numbers and the assumption (A5), the third term of (4.7)

$$\begin{aligned}
&-\frac{1}{n} \sum_{i=1}^n \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_{\alpha} \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} \right\} \right. \right. \\
&\quad \left. \left. \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \right] \right] \\
&\leq \frac{1}{n} \sum_{i=1}^n (e^{C_2 \|\mathbf{Y}_i\| + C_4} + C_5) \triangleq C_6
\end{aligned}$$

can be bounded by some constant  $C_6$  from above. Then (4.7) becomes

$$\begin{aligned}
0 &\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
&\quad + \frac{1}{n} \sum_{i=1}^n \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_{\alpha} \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} \right\} \right. \right. \\
&\quad \left. \left. \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^{V_i} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \right] \right] \\
&\quad + C_6 \\
&\leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
&\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_{\alpha} \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_{\alpha}^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_{\alpha} \right\} \right. \right. \\
&\quad \left. \left. \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^{\tau} e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} d\mathbf{b}_{\alpha 0} \right] \right]
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{n} \sum_{i=1}^n I(V_i \neq \tau) \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \left. \left. \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right\} \mathbf{b}_{\alpha 0} \right] \right] \\
& + C_6 \\
& \leq \frac{1}{n} \sum_{i=1}^n \sum_{s=1}^S \Delta_i I(S_i = s) \widehat{\zeta}_s \\
& + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\
& \quad \left. \left. \times \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \mathbf{b}_{\alpha 0} \right] \right] \\
& + C_7, \tag{4.8}
\end{aligned}$$

where  $C_7$  is a constant. On the other hand, since, for any  $\Gamma \geq 0$  and  $x > 0$ ,  $\Gamma \log(1 + x/\Gamma) \leq \Gamma x/\Gamma = x$ , we have that  $e^{-x} \leq (1 + x/\Gamma)^{-\Gamma}$ . Therefore, in the second term of (4.8),

$$\begin{aligned}
& \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\} \\
& \leq \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right\} \times \left\{ 1 + \frac{\sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t)}{\Gamma} \right\}^{-\Gamma} \\
& \leq \Gamma^\Gamma \times \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right\} \times \left\{ \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\}^{-\Gamma}. \tag{4.9}
\end{aligned}$$

Since  $Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}}) \geq -C_1 \|\mathbf{b}_{\alpha 0}\| - C_2 \|\mathbf{Y}_i\| - C_3$ ,

$$\begin{aligned}
\int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) & \geq \int_0^\tau e^{-C_1 \|\mathbf{b}_{\alpha 0}\| - C_2 \|\mathbf{Y}_i\| - C_3} d\widetilde{\Lambda}_s(t) \\
& = e^{-C_1 \|\mathbf{b}_{\alpha 0}\| - C_2 \|\mathbf{Y}_i\| - C_3} \times \{\widetilde{\Lambda}_s(\tau) - \widetilde{\Lambda}_s(0)\} \\
& = e^{-C_1 \|\mathbf{b}_{\alpha 0}\| - C_2 \|\mathbf{Y}_i\| - C_3}.
\end{aligned}$$

Thus, in (4.9),

$$\left\{ \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \int_0^\tau e^{Q_{1i\alpha}(t, \mathbf{b}_{\alpha 0}, \widehat{\boldsymbol{\theta}})} d\widetilde{\Lambda}_s(t) \right\}^{-\Gamma} \leq \left\{ \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \right\}^{-\Gamma} e^{C_1 \Gamma \|\mathbf{b}_{\alpha 0}\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma}$$

$$\text{and} \quad (4.9) \leq \Gamma^\Gamma \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \Gamma \log \left( \sum_{s=1}^S I(S_i = s) e^{\widehat{\zeta}_s} \right) \right. \\ \left. + C_1 \Gamma \|\mathbf{b}_{\alpha 0}\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma \right\}.$$

Therefore, (4.8) gives that

$$\begin{aligned} 0 &\leq C_7 + \frac{1}{n} \sum_{i=1}^n \Delta_i \left( \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right) \\ &\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \left[ \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\ &\quad \times \int_{\mathbf{b}_{\alpha 0}} \Gamma^\Gamma \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} - \Gamma \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right. \\ &\quad \left. \left. + C_1 \Gamma \|\mathbf{b}_{\alpha 0}\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma \right\} d\mathbf{b}_{\alpha 0} \right] \Bigg] \\ &= C_7 + \frac{1}{n} \sum_{i=1}^n \Delta_i \left( \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right) \\ &\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \log \left[ \Gamma^\Gamma \exp \left\{ -\Gamma \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right\} \sum_{\alpha=1}^K \left[ \widehat{w}_\alpha \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha \right\} \right. \right. \\ &\quad \times (2\pi)^{d_b/2} (2\pi)^{-d_b/2} \int_{\mathbf{b}_{\alpha 0}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha 0}^T \mathbf{b}_{\alpha 0} - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right. \\ &\quad \left. \left. + C_1 \Gamma \|\mathbf{b}_{\alpha 0}\| + C_2 \Gamma \|\mathbf{Y}_i\| + C_3 \Gamma \right\} d\mathbf{b}_{\alpha 0} \right] \Bigg] \\ &= C_7 + \frac{1}{n} \sum_{i=1}^n \Delta_i \left( \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right) + \frac{1}{2} \log(2\pi) \\ &\quad + \frac{1}{n} \sum_{i=1}^n I(V_i = \tau) \left[ \Gamma \log \Gamma - \Gamma \sum_{s=1}^S I(S_i = s) \widehat{\zeta}_s \right. \\ &\quad \left. + \log E_{\alpha, \mathbf{b}_0} \left[ \exp \left\{ \frac{1}{2} \mathbf{M}_{i\alpha}^T \mathbf{M}_{i\alpha} - \frac{1}{2} \widehat{\boldsymbol{\mu}}_\alpha^T \widehat{\boldsymbol{\Sigma}}_b^{-1} \widehat{\boldsymbol{\mu}}_\alpha - \sum_{j=1}^{n_i} \frac{B(\widehat{\boldsymbol{\beta}}; \mathbf{b}_{\alpha 0})}{A(D_i(t_j; \widehat{\phi}))} \right\} \right] \right] \end{aligned}$$

$$\begin{aligned}
& +C_1\Gamma\|\mathbf{b}_{\alpha 0}\|+C_2\Gamma\|\mathbf{Y}_i\|+C_3\Gamma\Big\}\Big]\Big] \\
= & C_8+\frac{1}{n}\sum_{i=1}^n\Delta_i\left(\sum_{s=1}^S\widehat{\zeta}_s\right)-\frac{\Gamma}{n}\sum_{i=1}^nI(V_i=\tau)\left(\sum_{s=1}^S\widehat{\zeta}_s\right)+C_9(\Gamma), \tag{4.10}
\end{aligned}$$

where  $C_8$  is a constant and  $C_9(\Gamma)$  is a deterministic function of  $\Gamma$ . For the  $s$ -th stratum, (4.10) is that

$$0 \leq C_8 + \sum_{i=1}^n \Delta_i I(S_i = s) \widehat{\zeta}_s - \frac{\Gamma}{n} \sum_{i=1}^n I(V_i = \tau) I(S_i = s) \widehat{\zeta}_s + C_9(\Gamma).$$

By the strong law of large numbers,  $\sum_{i=1}^n I(V_i = \tau) I(S_i = s)/n \rightarrow P(V_i = \tau, S_i = s) > 0$ . Then, we can choose  $\Gamma$  large enough such that  $\sum_{i=1}^n \Delta_i I(S_i = s)/n \leq (\Gamma/2n) \sum_{i=1}^n I(V_i = \tau) I(S_i = s)$ . Thus, we obtain that

$$0 \leq C_8 + C_9(\Gamma) - \frac{\Gamma}{2n} \sum_{i=1}^n I(V_i = \tau) I(S_i = s) \widehat{\zeta}_s.$$

In other words,

$$\widehat{\zeta}_s \leq \frac{(C_8 + C_9(\Gamma))2n}{\Gamma \sum_{i=1}^n I(V_i = \tau) I(S_i = s)} \rightarrow \frac{(C_8 + C_9(\Gamma))2}{\Gamma P(V_i = \tau, S_i = s)}.$$

If we denote  $B_{s0} = \exp\{2(C_8 + C_9(\Gamma))/(\Gamma P(V_i = \tau, S_i = s))\}$ , we conclude that  $\widehat{\Lambda}_s(\tau) \leq B_{s0}$ ,  $s = 1, \dots, S$ . Note that the above arguments hold for every sample in the probability space except a set with zero probability. Therefore, we have shown that, with probability one,  $\widehat{\Lambda}_s(\tau)$  is bounded for any sample size  $n$ .

In the third step, the goal of this step is to show that, if  $\widehat{\boldsymbol{\theta}} \rightarrow \boldsymbol{\theta}^*$  and  $\widehat{\Lambda}_s$  weakly converges to  $\Lambda_s^*$  with probability one, then  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$ ,  $s = 1, \dots, S$ . We set some preliminaries as the followings: For convenience, we omit the index  $i$  for subject and use  $\mathbf{O}$  to abbreviate the observed statistics  $(\mathbf{Y}, \mathbf{X}, \widetilde{\mathbf{X}}, V, \Delta, n_N, s)$  and  $\{\mathbf{Z}(t), \widetilde{\mathbf{Z}}(t), 0 \leq t \leq V\}$  for a subject. By dropping  $(\lambda_s(V))^\Delta$  from the complete data

likelihood function, we define

$$\begin{aligned}
& G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta} + \tilde{\mathbf{X}}_i \mathbf{b}_\alpha) - B(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} + C(Y_j; D(t_j; \phi)) \right] \right\} \\
&\times \exp \left\{ \Delta \left[ \tilde{\mathbf{Z}}(V)(\boldsymbol{\psi} \circ \mathbf{b}_\alpha) + \mathbf{Z}(V)\boldsymbol{\gamma} \right] - \int_0^V \exp \left\{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\boldsymbol{\gamma} \right\} d\Lambda_s(t) \right\} \\
&\times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_\alpha - \boldsymbol{\mu}_\alpha)^T \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_\alpha - \boldsymbol{\mu}_\alpha) \right\} w_\alpha, \\
\text{and } Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) &= \frac{\sum_\alpha \int_{\mathbf{b}_\alpha} G(\mathbf{b}_\alpha, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \exp \left\{ \tilde{\mathbf{Z}}(v)(\boldsymbol{\psi} \circ \mathbf{b}_\alpha) + \mathbf{Z}(v)\boldsymbol{\gamma} \right\} d\mathbf{b}_\alpha}{\sum_\alpha \int_{\mathbf{b}_\alpha} G(\mathbf{b}_\alpha, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b}_\alpha}.
\end{aligned}$$

Furthermore, for any measurable function  $f(\mathbf{O})$ , we use operator notation to define  $\mathbf{P}_n f = n^{-1} \sum_{i=1}^n f(\mathbf{O}_i)$  and  $\mathbf{P} f = \int f d\mathbf{P} = \mathbb{E}[f(\mathbf{O})]$ . Thus,  $\mathbf{P}_n f$  is the empirical measure from  $n$  i.i.d observations and  $\sqrt{n}(\mathbf{P}_n - \mathbf{P})$  is the empirical process based on these observations. We also define a class  $\mathcal{F} = \{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau], \boldsymbol{\theta} \in \Theta, \Lambda_s \in \mathbb{W}, \Lambda_s(0) = 0, \Lambda_s(\tau) \leq B_{s0}\}$ , where  $B_{s0}$  is the constant given in the second step and  $\mathbb{W}$  contains all nondecreasing functions in  $[0, \tau]$ . According to the result proved in Section 4.4.3.1,  $\mathcal{F}$  is P-Donsker.

Let  $m_s$  denote the number of subjects in stratum  $s$ ; i.e.  $n = \sum_{s=1}^S m_s$ .  $V_s$  and  $\Delta_s$  denote the observed time and censoring indicator for a subject belonging to stratum  $s$ , respectively. Thus,  $V_{sl}$  and  $\Delta_{sl}$  are the  $l$ -th subject observed time and censoring indicator in stratum  $s$ .

Now we start the proof of the third step. Since  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}})$  maximizes the function  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$ , where  $\boldsymbol{\Lambda} = (\Lambda_1, \dots, \Lambda_S)^T$  and  $\Lambda_s, s = 1, \dots, S$ , are any step functions with jumps only at  $V_i$  belonging to stratum  $s$  for which  $\Delta_i = 1$ , we differentiate  $l_n(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  with respect to  $\Lambda_s\{V_{sl}\}$  and obtain the following equation, satisfied by  $\widehat{\boldsymbol{\Lambda}}_s$ ,

$$\widehat{\boldsymbol{\Lambda}}_s\{V_{sl}\} = \frac{\Delta_{sl}}{m_s \mathbf{P}_{m_s} \left\{ I(V_s \geq v) Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s) \right\} \Big|_{v=V_{sl}}}$$

Imitating the above equation, we also can construct another function, denoted by  $\bar{\Lambda} = (\bar{\Lambda}_1, \dots, \bar{\Lambda}_S)^T$  such that  $\bar{\Lambda}_s$ ,  $s = 1, \dots, S$ , are also step functions with jumps only at the observed  $V_{sl}$  and the jump size  $\bar{\Lambda}_s\{V_{sl}\}$  is given by

$$\bar{\Lambda}_s\{V_{sl}\} = \frac{\Delta_{sl}}{m_s \mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_{sl}}}.$$

Equivalently,

$$\bar{\Lambda}_s(t) = \frac{1}{m_s} \sum_{l=1}^{m_s} \frac{I(V_{sl} \leq t) \Delta_{sl}}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_{sl}}}.$$

Then, we claim  $\bar{\Lambda}_s(t)$  uniformly converges to  $\Lambda_{s0}(t)$  in  $[0, \tau]$ . To prove the claim, note that

$$\begin{aligned} & \sup_{t \in [0, \tau]} \left| \bar{\Lambda}_s(t) - \mathbf{E} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] \right| \\ &= \sup_{t \in [0, \tau]} \left| \frac{1}{m_s} \sum_{l=1}^{m_s} \frac{I(V_{sl} \leq t) \Delta_{sl}}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_{sl}}} \right. \\ & \quad \left. - \mathbf{P}_{m_s} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] \right. \\ & \quad \left. + \mathbf{P}_{m_s} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] - \mathbf{P} \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] \right| \\ &\leq \sup_{t \in [0, \tau]} \left| \frac{1}{m_s} \sum_{l=1}^{m_s} I(V_{sl} \leq t) \Delta_{sl} \left[ \frac{1}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right. \right. \\ & \quad \left. \left. - \frac{1}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_{sl}}} \right] \right| \\ & \quad + \sup_{t \in [0, \tau]} \left| \left( \mathbf{P}_{m_s} - \mathbf{P} \right) \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] \right| \\ &\leq \sup_{t \in [0, \tau]} \left| \frac{1}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} - \frac{1}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}} \right| \\ & \quad + \sup_{t \in [0, \tau]} \left| \left( \mathbf{P}_{m_s} - \mathbf{P} \right) \left[ \frac{I(V_s \leq t) \Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \big|_{v=V_s}} \right] \right|. \tag{4.11} \end{aligned}$$

In (4.11), the right hand side converges to 0 because the first and second terms on the right hand side converges to 0 in the following: First, according to Section 4.4.3.1,  $\{Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) : v \in [0, \tau]\}$  is a bounded and Glivenko-Cantelli class.  $\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) : v \in [0, \tau]\}$  is also a Glivenko-Cantelli class because  $\{I(V_s \geq v) : v \in [0, \tau]\}$  is a Glivenko-Cantelli class and the functional  $(f, g) \rightarrow fg$  for any bounded two functions  $f$  and  $g$  is Lipschitz continuous. Then, we obtain that  $\sup_{t \in [0, \tau]} \left| \mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} - \mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} \right|$  converges to 0. Besides, from Section 4.4.3.1,  $\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\} > \mathbf{P} \{I(V_s \geq v) \exp\{-C_{10} - C_{11}\|\mathbf{Y}\|\}\}$  for the two constants  $C_{10}$  and  $C_{11}$ , which means  $\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}$  is bounded from below. Thus, the first term tends to 0. Second, since the class  $\{I(V_s \leq t)\Delta_s / \mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}|_{v=V_s} : t \in [0, \tau]\}$  is also a Glivenko-Cantelli class, the second term vanishes as  $m_s$  goes to infinity.

Therefore, we conclude that  $\bar{\Lambda}_s(t)$  uniformly converges to

$$\mathbf{E} \left[ \frac{I(V_s \leq t)\Delta_s}{\mathbf{P} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}|_{v=V_s}} \right]. \quad (4.12)$$

We can easily verify that (4.12) is equal to  $\Lambda_{s0}(t)$ . Thus, the claim that  $\bar{\Lambda}_s(t)$  uniformly converges to  $\Lambda_{s0}(t)$  in  $[0, \tau]$  has been proved.

From the construction of  $\bar{\Lambda}_s(t)$ , we obtain

$$\widehat{\Lambda}_s(t) = \int_0^t \frac{d\widehat{\Lambda}_s(v)}{d\bar{\Lambda}_s(v)} d\bar{\Lambda}_s(v) = \int_0^t \frac{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}_{m_s} \{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}} d\bar{\Lambda}_s(v). \quad (4.13)$$

$\widehat{\Lambda}_s(t)$  is absolutely continuous with respect to  $\bar{\Lambda}_s(t)$ . On the other hand, since both  $\{I(V_s \geq v) : v \in [0, \tau]\}$  and  $\mathcal{F}$  are Glivenko-Cantelli classes,  $\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau]\}$  is also a Glivenko-Cantelli class. Thus, we have

$$\begin{aligned} & \sup_{v \in [0, \tau]} |(\mathbf{P}_{m_s} - \mathbf{P})\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}| + \sup_{v \in [0, \tau]} |(\mathbf{P}_{m_s} - \mathbf{P})\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}| \\ & \longrightarrow 0 \quad \text{a.s.} \end{aligned}$$

By the bounded convergence theorem and the fact that  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}^*$  and  $\widehat{\Lambda}_s$  converges to  $\Lambda_s^*$ , for each  $v$ ,  $\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\} \longrightarrow \mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}$ ; moreover, it is straightforward to check the derivative of  $\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}$  with respect to  $v$ . Thus, by the Arzela-Ascoli theorem, uniformly in  $[0, \tau]$ ,

$$\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\} \longrightarrow \mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}.$$

Then, combining the above result and (4.13), it holds that, uniformly in  $[0, \tau]$ ,

$$\frac{d\widehat{\Lambda}_s(v)}{d\Lambda_s(v)} = \frac{\mathbf{P}_{m_s}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}_{m_s}\{I(V_s \geq v)Q(v, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)\}} \longrightarrow \frac{\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}}. \quad (4.14)$$

After taking limits on both sides of (4.13), we obtain that

$$\lim_{m_s \rightarrow \infty} \widehat{\Lambda}_s(t) = \int_0^t \frac{\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}}{\mathbf{P}\{I(V_s \geq v)Q(v, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)\}} d\Lambda_{s0}(v), \quad (4.15)$$

Therefore, since  $\Lambda_{s0}(t)$  is differentiable with respect to the Lebesgue measure, so is  $\Lambda_s^*(t)$ ; that is, (4.15) is equal to

$$\int_0^t \frac{d\Lambda_s^*(v)}{d\Lambda_{s0}(v)} d\Lambda_{s0}(v). \quad (4.16)$$

And we denote  $\lambda_s^*(t)$  as the derivative of  $\Lambda_s^*(t)$ . Additionally, from (4.14) ~ (4.16), note that  $\widehat{\Lambda}_s\{V_s\}/\widehat{\Lambda}_s\{V_s\}$  uniformly converges to  $d\Lambda_s^*(V_s)/d\Lambda_{s0}(V_s) = \lambda_s^*(V_s)/\lambda_{s0}(V_s)$ . Therefore, a second conclusion is that  $\widehat{\Lambda}_s$  uniformly converges to  $\Lambda_s^*$  since  $\Lambda_s^*$  is contin-



uous.

On the other hand,

$$\begin{aligned}
& n^{-1}l_n(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}) - n^{-1}l_n(\boldsymbol{\theta}_0, \bar{\boldsymbol{\Lambda}}) \\
&= \sum_{s=1}^S \left( \mathbf{P}_{m_s} \left[ \Delta_s \log \frac{\widehat{\Lambda}_s\{V_s\}}{\bar{\Lambda}_s\{V_s\}} \right] + \mathbf{P}_{m_s} \left[ \log \frac{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}_{\alpha}}{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}_{\alpha}} \right] \right) \\
&\geq 0.
\end{aligned} \tag{4.17}$$

Using the result of Section 4.4.3.1 and similar arguments as above, we can verify that

$$\log \frac{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}_{\alpha}}{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}_{\alpha}}$$

belongs to a Glivenko-Cantelli class and

$$\mathbf{P} \left[ \log \frac{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s) d\mathbf{b}_{\alpha}}{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \bar{\Lambda}_s) d\mathbf{b}_{\alpha}} \right] \rightarrow \mathbf{P} \left[ \log \frac{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha}}{\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha}} \right].$$

Since  $\widehat{\Lambda}_s\{V_s\}/\bar{\Lambda}_s\{V_s\}$  uniformly converges to  $\lambda_s^*(V_s)/\lambda_{s0}(V_s)$ , we obtain that, from (4.17),

$$\mathbf{P} \left[ \log \left\{ \frac{(\lambda_s^*(V_s))^{\Delta_s} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha}}{(\lambda_{s0}(V_s))^{\Delta_s} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha}} \right\} \right] \geq 0.$$

Note that the left-hand side of the inequality is the negative Kullback-Leibler information. Then, the equality holds with probability one, and it immediately follows

$$(\lambda_s^*(V_s))^{\Delta_s} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} = (\lambda_{s0}(V_s))^{\Delta_s} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha}. \tag{4.18}$$

Our proof will be completed if we can show  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_0$  and  $\Lambda_s^* = \Lambda_{s0}$  from (4.18). Since (4.18) holds with probability one, (4.18) holds for any  $(V_s, \Delta_s = 1)$  and the case  $(V_s = \tau, \Delta_s = 0)$ , but may not hold for  $(V_s, \Delta_s = 0)$  when  $V_s \in (0, \tau)$ . However, we can show

that (4.18) is also true for  $(V_s, \Delta_s = 0)$  when  $V_s \in (0, \tau)$ . To do this, treating both sides of (4.18) as functions of  $V_s$ , we integrate these functions over an interval  $(V_s, \tau)$  for  $\Delta_s = 0$  as the following;

$$\int_{V_s}^{\tau} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} = \int_{V_s}^{\tau} \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha}$$

to obtain that

$$\begin{aligned} & \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=\tau} - \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=V_s} \\ &= \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=\tau} - \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=V_s}. \end{aligned}$$

After comparing this above equality with another following equality, which is given by (4.18) at  $\Delta_s = 0$  and  $V_s = \tau$ ,

$$\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=\tau} = \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=\tau},$$

we obtain

$$\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=V_s} = \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=V_s},$$

and therefore

$$\sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0} = \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0};$$

that is, (4.18) also holds for any  $V_s$  and  $\Delta_s = 0$ .

Thus, first to show that  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_0$ ,  $\boldsymbol{\phi}^* = \boldsymbol{\phi}_0$ ,  $\boldsymbol{\Sigma}_b^* = \boldsymbol{\Sigma}_{b0}$ ,  $\boldsymbol{\mu}_{\alpha}^* = \boldsymbol{\mu}_{\alpha 0}$ , and  $w_{\alpha}^* = w_{\alpha 0}$ ,  $\alpha = 1, \dots, K$ , we let  $\Delta_s = 0$  and  $V_s = 0$  in (4.18). After integrating over  $\mathbf{b}_{\alpha}$  and summing

up over  $\alpha$ , we have that, with probability one,

$$\begin{aligned}
& \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=0} = \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\mathbf{b}_{\alpha} \Big|_{\Delta_s=0, V_s=0} \\
\Rightarrow & \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta}^* + \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}) - B(\boldsymbol{\beta}^*, \mathbf{b}_{\alpha})}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} \\
& \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b^*|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha}^*)^T \boldsymbol{\Sigma}_b^{*-1} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha}^*) \right\} d\mathbf{b}_{\alpha} \\
& = \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}) - B(\boldsymbol{\beta}_0, \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} + C(Y_j; D(t_j; \phi_0)) \right] \right\} \\
& \quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_{b0}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha 0})^T \boldsymbol{\Sigma}_{b0}^{-1} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha 0}) \right\} d\mathbf{b}_{\alpha} \\
\Rightarrow & \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b^*|^{-1/2} \\
& \quad \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}}{A(D(t_j; \phi^*))} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_{\alpha})}{A(D(t_j; \phi^*))} \right. \\
& \quad \left. - \frac{1}{2} \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_b^{*-1} \mathbf{b}_{\alpha} + \boldsymbol{\mu}_{\alpha}^{*T} \boldsymbol{\Sigma}_b^{*-1} \mathbf{b}_{\alpha} - \frac{1}{2} \boldsymbol{\mu}_{\alpha}^{*T} \boldsymbol{\Sigma}_b^{*-1} \boldsymbol{\mu}_{\alpha}^* \right\} w_{\alpha}^* d\mathbf{b}_{\alpha} \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}_0}{A(D(t_j; \phi_0))} + C(Y_j; D(t_j; \phi_0)) \right] \right\} (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_{b0}|^{-1/2} \\
& \quad \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}}{A(D(t_j; \phi_0))} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right. \\
& \quad \left. - \frac{1}{2} \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_{\alpha} + \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_{\alpha} - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha 0} \right\} w_{\alpha 0} d\mathbf{b}_{\alpha}.
\end{aligned}$$

By some algebra, the left hand side becomes

$$\begin{aligned}
& \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b^*|^{-1/2} \\
& \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \left[ (\boldsymbol{\Sigma}_b^{*-1/2} \mathbf{b}_{\alpha})^T (\boldsymbol{\Sigma}_b^{*-1/2} \mathbf{b}_{\alpha}) \right. \right. \\
& \quad \left. \left. - 2 \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_{\alpha}^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \boldsymbol{\Sigma}_b^{*-1/2} \mathbf{b}_{\alpha} \right. \right. \\
& \quad \left. \left. + \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_{\alpha}^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right] \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_{\alpha}^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right]^T \right] \right\} d\mathbf{b}_{\alpha}
\end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2} \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right] \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right]^T \\
& - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} - \frac{1}{2} \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \boldsymbol{\mu}_\alpha^* \Big\} w_\alpha^* d\mathbf{b}_\alpha \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b^*|^{-1/2} \\
& \times \sum_{\alpha} w_\alpha^* \left[ \exp \left\{ \frac{1}{2} \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right] \right. \right. \\
& \quad \times \left[ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^{*1/2} \right]^T - \frac{1}{2} \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \boldsymbol{\mu}_\alpha^* \Big\} \\
& \quad \times \int_{\mathbf{b}_\alpha} \exp \left\{ -\frac{1}{2} \left[ \boldsymbol{\Sigma}_b^{*-1/2} \mathbf{b}_\alpha - \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right)^T \right]^T \right. \\
& \quad \quad \times \left. \left[ \boldsymbol{\Sigma}_b^{*-1/2} \mathbf{b}_\alpha - \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right)^T \right] \right\} \\
& \quad \times \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} d\mathbf{b}_\alpha \Big] \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} \\
& \quad \times \sum_{\alpha} w_\alpha^* \left[ \exp \left\{ \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} + \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right)^T \right. \right. \\
& \quad \quad \left. \left. - \frac{1}{2} \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \boldsymbol{\mu}_\alpha^* \right\} \mathbb{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \left( \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right) \right\} \right] \right] \\
& = \exp \left\{ \sum_{j=1}^{n_N} \left[ \frac{Y_j \mathbf{X}_j \boldsymbol{\beta}^*}{A(D(t_j; \phi^*))} + C(Y_j; D(t_j; \phi^*)) \right] \right\} \\
& \quad \times \sum_{\alpha} w_\alpha^* \left[ \exp \left\{ \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right. \right. \\
& \quad \quad \left. \left. + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \boldsymbol{\Sigma}_b^* \left( \boldsymbol{\mu}_\alpha^{*T} \boldsymbol{\Sigma}_b^{*-1} \right)^T \right\} \times \mathbb{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} \right] \right] \\
& = \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi^*)) \right\} \sum_{\alpha} w_\alpha^* \left[ \exp \left\{ \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right. \right. \\
& \quad \quad \left. \left. + \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi^*))} \left( \mathbf{X}_j \boldsymbol{\beta}^* + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_\alpha^* \right) \right\} \times \mathbb{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} \right] \right] \\
& = \sum_{\alpha} \left[ \exp \left\{ \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T \right. \right.
\end{aligned}$$

$$\begin{aligned}
& + \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi^*))} \left( \mathbf{X}_j \boldsymbol{\beta}^* + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_\alpha^* \right) \Bigg\} \\
& \times \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi^*)) \right\} w_\alpha^* \mathbf{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} \right] \Bigg]. \quad (4.19)
\end{aligned}$$

Likewise, the right-hand side becomes

$$\begin{aligned}
& \sum_\alpha \left[ \exp \left\{ \frac{1}{2} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \boldsymbol{\Sigma}_{b0} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T \right. \right. \\
& \quad \left. \left. + \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \left( \mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right) \right\} \right. \\
& \quad \left. \times \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi_0)) \right\} w_{\alpha 0} \mathbf{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right] \right]. \quad (4.20)
\end{aligned}$$

Then, to compare the coefficients of  $\mathbf{Y}^T \mathbf{Y}$  and  $\mathbf{Y}$  in the exponential part and the constant term out of the exponential part from (4.19) and (4.20), we have

$$\left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi^*))} \right)^T = \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \boldsymbol{\Sigma}_{b0} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T, \quad (4.21)$$

$$\sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi^*))} \left( \mathbf{X}_j \boldsymbol{\beta}^* + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_\alpha^* \right) = \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \left( \mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right), \quad (4.22)$$

and

$$\begin{aligned}
& \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi^*)) \right\} w_\alpha^* \mathbf{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} \right] \\
& = \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi_0)) \right\} w_{\alpha 0} \mathbf{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right]. \quad (4.23)
\end{aligned}$$

Furthermore, by the assumption of the generalized linear mixed model with canonical link function for longitudinal outcome  $Y(t)$  at time  $t$ , we have  $\mu(t) = \mathbf{E}(Y(t) | \mathbf{b}) = B'(\eta(t))$  and  $v(t) = \text{Var}(Y(t) | \mathbf{b}) = B''(\eta(t))A(\phi(t))$ , where  $\mathbf{b} = \sum_{k=1}^K I(\alpha = k) \mathbf{b}_k$ ,  $\eta(t) = g(\mu(t)) = \mathbf{X}(t) \boldsymbol{\beta} + \tilde{\mathbf{X}}(t) \mathbf{b}$ ,  $v(t) = v(\mu(t))A(\phi(t))$ ,  $g(\cdot)$  and  $v(\cdot)$  are known link and

variance functions respectively, and  $B'(\eta(t))$  and  $B''(\eta(t))$  are the first and second derivatives of  $B(\eta(t))$  with respect to the canonical parameter  $\eta(t)$ . Hence, we have

$$E(Y_j|\mathbf{b}) = B'(\eta_j) = B'(\boldsymbol{\beta}^*; \mathbf{b}) = B'(\boldsymbol{\beta}_0; \mathbf{b}) \quad (4.24)$$

and

$$\text{Var}(Y_j|\mathbf{b}) = B''(\eta_j)A(D(t_j; \phi)) = B''(\boldsymbol{\beta}^*; \mathbf{b})A(D(t_j; \phi^*)) = B''(\boldsymbol{\beta}_0; \mathbf{b})A(D(t_j; \phi_0)). \quad (4.25)$$

By the continuous mapping theorem and (4.24), we obtain  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_0$ . Then, (4.25) becomes  $B''(\boldsymbol{\beta}_0; \mathbf{b})A(D(t_j; \phi^*)) = B''(\boldsymbol{\beta}_0; \mathbf{b})A(D(t_j; \phi_0))$ . Hence, by assumption (A6),  $A(D(t_j; \phi^*)) = A(D(t_j; \phi_0))$ , and, by the continuous mapping theorem, we obtain  $D(t_j; \phi^*) = D(t_j; \phi_0)$ ,  $j = 1, \dots, n_N$ , and  $\phi^* = \phi_0$ . Thus, (4.21) can be written as

$$\left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T = \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right) \boldsymbol{\Sigma}_{b0} \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi_0))} \right)^T.$$

Then, by assumption (A6), we obtain  $\boldsymbol{\Sigma}_b^* = \boldsymbol{\Sigma}_{b0}$ . Since  $\boldsymbol{\beta}^* = \boldsymbol{\beta}_0$  and  $\phi^* = \phi_0$ , (4.22) can be written as

$$\sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \left( \mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_\alpha^* \right) = \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \left( \mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right).$$

Also, by assumption (A6), we obtain  $\boldsymbol{\mu}_\alpha^* = \boldsymbol{\mu}_{\alpha 0}$ ,  $\alpha = 1, \dots, K$ . In (4.23) for the constant terms, note that the random effect  $\mathbf{b}_\alpha$  on the left-hand side follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_b^* \left( \sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi^*)) \right)^T + \boldsymbol{\mu}_\alpha^*$  and covariance  $\boldsymbol{\Sigma}_b^*$  and the random effect  $\mathbf{b}_\alpha$  on the right-hand side follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_{b0} \left( \sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0)) \right)^T + \boldsymbol{\mu}_{\alpha 0}$  and covariance  $\boldsymbol{\Sigma}_{b0}$ . (i) Because  $\boldsymbol{\Sigma}_b^* = \boldsymbol{\Sigma}_{b0}$ ,  $\boldsymbol{\mu}_\alpha^* = \boldsymbol{\mu}_{\alpha 0}$ ,  $\alpha = 1, \dots, K$ , and  $\phi^* = \phi_0$ , the random effects  $\mathbf{b}_\alpha$ 's on both sides

follow the same multivariate normal distribution. (ii) Besides, because  $\beta^* = \beta_0$  and  $\phi^* = \phi_0$ , we have  $\sum_{j=1}^{n_N} \frac{B(\beta^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} = \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))}$ . By (i) and (ii), we obtain (iii)  $E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ \frac{B(\beta^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} \right\} \right] = E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right]$ . Also, (iv) since  $\phi^* = \phi_0$ , we have  $\exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi^*)) \right\} = \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi_0)) \right\}$ . By (iii) and (iv), (4.23) can be written as

$$\begin{aligned} & \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi_0)) \right\} w_\alpha^* E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right] \\ &= \exp \left\{ \sum_{j=1}^{n_N} C(Y_j; D(t_j; \phi_0)) \right\} w_{\alpha 0} E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right]. \end{aligned}$$

Then, by assumption (A6), we obtain  $w_\alpha^* = w_{\alpha 0}$ .  $\alpha = 1, \dots, K$ .

Next, to show that  $\psi^* = \psi_0$ ,  $\gamma^* = \gamma_0$  and  $\Lambda_s^* = \Lambda_{s0}$ , we let  $\Delta_s = 0$  in (4.18). Through the similar arguments done for the proof of  $\beta^* = \beta_0$ ,  $\phi^* = \phi_0$ ,  $\Sigma_b^* = \Sigma_{b0}$ ,  $\mu_\alpha^* = \mu_{\alpha 0}$ , and  $w_\alpha^* = w_{\alpha 0}$ ,  $\alpha = 1, \dots, K$ , we obtain

$$\begin{aligned} & E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma^* \} d\Lambda_s^*(t) \right\} \right] \\ &= E_{\mathbf{b}_\alpha|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma_0 \} d\Lambda_{s0}(t) \right\} \right], \end{aligned} \tag{4.26}$$

where the random effects  $\mathbf{b}_\alpha$ 's on both sides follow a multivariate normal distribution with mean  $\Sigma_{b0} \left( \sum_{j=1}^{n_N} Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0)) \right)^T + \mu_{\alpha 0}$  and covariance  $\Sigma_{b0}$ .

For any fixed  $\tilde{\mathbf{X}} = (\tilde{\mathbf{X}}_1^T, \dots, \tilde{\mathbf{X}}_{n_N}^T)^T$ , treating  $\tilde{\mathbf{X}}^T \mathbf{Y}$  as a parameter in this normal family,  $\mathbf{b}_\alpha = \sum_{k=1}^K I(\alpha = k) \mathbf{b}_k$  is the complete statistic for  $\tilde{\mathbf{X}}^T \mathbf{Y}$ . Therefore,

$$\begin{aligned} & \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta^*; \mathbf{b}_\alpha)}{A(D(t_j; \phi^*))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma^* \} d\Lambda_s^*(t) \right\} \\ &= \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\beta_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma_0 \} d\Lambda_{s0}(t) \right\}. \end{aligned}$$

Since  $\beta^* = \beta_0$  and  $\phi^* = \phi_0$ , equivalently, we have

$$\exp \left\{ \tilde{\mathbf{Z}}(t)(\psi^* \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma^* \right\} \lambda_s^*(t) = \exp \left\{ \tilde{\mathbf{Z}}(t)(\psi_0 \circ \mathbf{b}_\alpha) + \mathbf{Z}(t)\gamma_0 \right\} \lambda_{s0}(t).$$

By assumptions (A3) and (A6),  $\psi^* = \psi_0$ ,  $\gamma^* = \gamma_0$  and  $\Lambda_s^* = \Lambda_{s0}$ .

Since all the three steps are completed, we can conclude that, with probability one,  $\widehat{\boldsymbol{\theta}}$  converges to  $\boldsymbol{\theta}_0$  and  $\widehat{\boldsymbol{\Lambda}}$  converges to  $\boldsymbol{\Lambda}_0$  in  $[0, \tau]$ . Moreover, as mentioned in the beginning of this proof for consistency, since  $\boldsymbol{\Lambda}_0$  is continuous in  $[0, \tau]$ , the latter can be strengthened to uniform convergence; that is,  $\sup_{t \in [0, \tau]} \|\widehat{\boldsymbol{\Lambda}}(t) - \boldsymbol{\Lambda}_0(t)\| \rightarrow 0$  almost surely. Therefore, Theorem 4.1 is proved.

#### 4.4.2 Proof of asymptotic normality

Asymptotic distribution for the proposed estimator can be shown if we can verify the conditions of Theorem 3.3.1 (p310) in van der Vaart and Wellner (1996). Then, we will show that the distribution is normal. For completeness, we state this theorem below following Theorem 4 in Appendix A of Parner (1998).

**Theorem 4.3.** (Theorem 3.3.1 in van der Vaart and Wellner, 1996; Theorem 4 in Parner, 1998) *Let  $U_n$  and  $U$  be random maps and a fixed map, respectively, from  $\xi$  to a Banach space such that:*

- (a)  $\sqrt{n}(U_n - U)(\widehat{\xi}_n) - \sqrt{n}(U_n - U)(\xi_0) = o_P^*(1 + \sqrt{n}\|\widehat{\xi}_n - \xi_0\|)$ .
- (b) *The sequence  $\sqrt{n}(U_n - U)(\xi_0)$  converges in distribution to a tight random element  $\mathbf{Z}$ .*
- (c) *the function  $\xi \rightarrow U(\xi)$  is Fréchet differentiable at  $\xi_0$  with a continuously invertible derivative  $\nabla U_{\xi_0}$  (on its range).*
- (d)  $U_{\xi_0}$  and  $\widehat{\xi}_n$  satisfies  $U_n(\widehat{\xi}_n) = o_P^*(n^{-1/2})$  and converges in outer probability to  $\xi_0$ .



Then  $\sqrt{n}(\widehat{\xi}_n - \xi_0) \Rightarrow \nabla U_{\xi_0}^{-1} \mathbf{Z}$ .

We will prove the conditions (a)~(d). In our situation, the parameter  $\xi_s = (\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) \in \Xi = \{(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) : \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, s = 1, \dots, S\}$  for a fixed small constant  $\delta$ . We note that  $\Xi$  is a convex set. Define a set  $\mathcal{H} = \{(\mathbf{h}_1, h_2) : \|\mathbf{h}_1\| \leq 1, \|h_2\|_V \leq 1\}$ , where  $\|h_2\|_V$  is the total variation of  $h_2$  in  $[0, \tau]$  defined as

$$\sup_{0=t_0 \leq t_1 \leq \dots \leq t_l = \tau} \sum_{j=1}^l |h_2(t_j) - h_2(t_{j-1})|.$$

Furthermore, we define that, for stratum  $s$ ,

$$U_{m_s}(\xi_s)(\mathbf{h}_1, h_2) = \mathbf{P}_{m_s} \{l_\theta(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)^T \mathbf{h}_1 + l_{\boldsymbol{\Lambda}_s}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)[h_2]\}$$

and

$$U_s(\xi_s)(\mathbf{h}_1, h_2) = \mathbf{P} \{l_\theta(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)^T \mathbf{h}_1 + l_{\boldsymbol{\Lambda}_s}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)[h_2]\},$$

where  $l_\theta(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)$  is the first derivative of the log-likelihood function from one single subject belonging to stratum  $s$ , denoted by  $l(\mathbf{O}; \boldsymbol{\theta}, \boldsymbol{\Lambda}_s)$ , with respect to  $\boldsymbol{\theta}$ , and  $l_{\boldsymbol{\Lambda}_s}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)$  is the derivative of  $l(\mathbf{O}; \boldsymbol{\theta}, \boldsymbol{\Lambda}_{s\varepsilon})$  at  $\varepsilon = 0$ , where  $\boldsymbol{\Lambda}_{s\varepsilon}(t) = \int_0^t (1 + \varepsilon h_2(u)) d\boldsymbol{\Lambda}_s(u)$ . Therefore, we can see that both  $U_{m_s}$  and  $U_s$  map from  $\Xi$  to  $\ell^\infty(\mathcal{H})$  and  $\sqrt{m_s}\{U_{m_s}(\xi_s) - U_s(\xi_s)\}$  is an empirical process in the space  $\ell^\infty(\mathcal{H})$ .

Denote  $(\mathbf{h}_1^\beta, \mathbf{h}_1^\phi, \mathbf{h}_1^{\Sigma_b}, \mathbf{h}_1^\mu, \mathbf{h}_1^w, \mathbf{h}_1^\psi, \mathbf{h}_1^\gamma)$  as the corresponding components of  $\mathbf{h}_1$  for the parameters  $(\boldsymbol{\beta}, \boldsymbol{\phi}, \text{Vec}(\boldsymbol{\Sigma}_b), \boldsymbol{\mu}, \mathbf{w}, \boldsymbol{\psi}, \boldsymbol{\gamma})$ , respectively. From Section 4.4.3.2, for any  $(\mathbf{h}_1, h_2) \in \mathcal{H}$ , the class

$$\begin{aligned} \mathcal{G} &= \{l_\theta(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)^T \mathbf{h}_1 + l_{\boldsymbol{\Lambda}_s}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \boldsymbol{\Lambda}_{s0})^T \mathbf{h}_1 + l_{\boldsymbol{\Lambda}_s}(\boldsymbol{\theta}_0, \boldsymbol{\Lambda}_{s0})[h_2], \\ &\quad \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, (\mathbf{h}_1, h_2) \in \mathcal{H}\} \end{aligned}$$

is shown as P-Donsker (Section 2.1 of van der Vaart and Wellner (1996), and it is also

implied that

$$\sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \mathbf{P} \left[ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right]^2 \longrightarrow 0$$

as  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \rightarrow 0$ . Then we conclude the followings:

- (a) follows from Lemma 3.3.5 (p311) of van der Vaart and Wellner (1996).
- (b) holds as a result of Section 4.4.3.2 and the convergence is defined in the metric space  $\ell^\infty(\mathcal{H})$  by the Donsker theorem (Section 2.5 of van der Vaart and Wellner (1996)).
- (d) is true because  $(\widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)$  maximizes  $\mathbf{P}_{m_s} l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$ ,  $(\boldsymbol{\theta}_0, \Lambda_{s0})$  maximizes  $\mathbf{P} l(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$ , and  $(\widehat{\boldsymbol{\theta}}, \widehat{\Lambda}_s)$  converges to  $(\boldsymbol{\theta}_0, \Lambda_{s0})$  from Theorem 4.1.

Now, we need to verify the conditions in (c). Since the proof of the first half in (c), that the function  $\xi \rightarrow U(\xi)$  is *Fréchet* differentiable at  $\xi_0$ , is given in Section 4.4.3.3, we will only prove that the derivative  $\nabla U_{\xi_0}$  is continuously invertible on its range  $\ell^\infty(\mathcal{H})$ . According to Section 4.4.3.3,  $\nabla U_{\xi_0}$  can be expressed as follows: for any  $(\boldsymbol{\theta}_1, \Lambda_{s1})$  and  $(\boldsymbol{\theta}_2, \Lambda_{s2})$  in  $\Xi$ ,

$$\nabla U_{\xi_0}(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\Lambda_{s1} - \Lambda_{s2})(t), \quad (4.27)$$

where both  $\Omega_1$  and  $\Omega_2$  are linear operators on  $\mathcal{H}$ , and  $\Omega = (\Omega_1, \Omega_2)$  maps  $\mathcal{H} \subset \mathbf{R}^d \times \text{BV}[0, \tau]$  to  $\mathbf{R}^d \times \text{BV}[0, \tau]$ , where  $\text{BV}[0, \tau]$  contains all the functions with finite total variation in  $[0, \tau]$ . The explicit expressions of  $\Omega_1$  and  $\Omega_2$  are given in Section 4.4.3.3. From (4.27), we can treat  $(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})$  as an element in  $\ell^\infty(\mathcal{H})$  via the following definition:

$$(\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2, \Lambda_{s1} - \Lambda_{s2})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\Lambda_{s1} - \Lambda_{s2})(t), \quad \forall (\mathbf{h}_1, h_2) \in \mathbf{R}^d \times \text{BV}[0, \tau].$$

Then  $\nabla U_{\xi_0}$  can be expanded as a linear operator from  $\ell^\infty(\mathcal{H})$  to itself. Therefore, if we can show that there exists some positive constant  $\varepsilon$  such that  $\varepsilon\mathcal{H} \subset \Omega(\mathcal{H})$ , then we will have that for any  $(\delta\boldsymbol{\theta}, \delta\Lambda_s) \in \ell^\infty(\mathcal{H})$ ,

$$\begin{aligned} \|\nabla U_{\xi_0}(\delta\boldsymbol{\theta}, \delta\Lambda_s)\|_{\ell^\infty(\mathcal{H})} &= \sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \left| \delta\boldsymbol{\theta}^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d\delta\Lambda_s(t) \right| \\ &= \|(\delta\boldsymbol{\theta}, \delta\Lambda_s)\|_{\ell^\infty(\Omega(\mathcal{H}))} \geq \varepsilon \|(\delta\boldsymbol{\theta}, \delta\Lambda_s)\|_{\ell^\infty(\mathcal{H})}, \end{aligned}$$

and  $\nabla U_{\xi_0}$  will be continuously invertible.

Note that to prove  $\varepsilon\mathcal{H} \subset \Omega(\mathcal{H})$  for some  $\varepsilon$  is equivalent to showing that  $\Omega$  is invertible. We also note from Section 4.4.3.3, that  $\Omega$  is the summation of an invertible operator and a compact operator. By Theorem 4.25 of Rudin (1973), for the proof of the invertibility of  $\Omega$ , it is sufficient to verify that  $\Omega$  is one to one: if  $\Omega[\mathbf{h}_1, h_2] = 0$ , then, by choosing  $\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2 = \varepsilon^* \mathbf{h}_1$  and  $\Lambda_{s1} - \Lambda_{s2} = \varepsilon^* \int h_2 d\Lambda_{s0}$  in (4.27) for a small constant  $\varepsilon^*$ , we obtain

$$\nabla U_{\xi_0}(\mathbf{h}_1, \int h_2 d\Lambda_{s0})[\mathbf{h}_1, h_2] = \varepsilon^* (\mathbf{h}_1^T, h_2) \begin{pmatrix} \Omega_1[\mathbf{h}_1, h_2] \\ \Omega_2[\mathbf{h}_1, h_2] \end{pmatrix} = \varepsilon^* (\mathbf{h}_1^T, h_2) \Omega[\mathbf{h}_1, h_2] = 0.$$

By the definition of  $\nabla U_{\xi_0}$ , we note that  $\nabla U_{\xi_0}(\mathbf{h}_1, \int h_2 d\Lambda_{s0})[\mathbf{h}_1, h_2]$  is the negative information matrix in the submodel  $(\boldsymbol{\theta}_0 + \varepsilon \mathbf{h}_1, \Lambda_{s0} + \varepsilon \int h_2 d\Lambda_{s0})$ . Thus, the score function along this submodel should be zero with probability one; that is,  $l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] = 0$ ; that is, with probability one, for the numerator of the score function

$$\begin{aligned} 0 &= \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \sum_{j=1}^{n_N} \frac{1}{A(D(t_j; \phi_0))} (Y_j \mathbf{X}_j - B'(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})) \mathbf{h}_1^{\beta} \right. \\ &\quad \left. + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}) - B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))^2} \right) A'(D(t_j; \phi_0)) + C'(Y_j; D(t_j; \phi_0)) \right\} \mathbf{h}_1^{\phi} \right] \end{aligned}$$

$$\begin{aligned}
& + \frac{1}{2}(\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0})^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} (\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0}) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \\
& + \left( \mathbf{b}_\alpha - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0} \right)^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_{\alpha 0}} + \Delta_s \{ (\tilde{\mathbf{Z}}(V_s) \circ \mathbf{b}_\alpha^T) \mathbf{h}_1^\psi + \mathbf{Z}(V_s) \mathbf{h}_1^\gamma \} \\
& - \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t) (\boldsymbol{\psi}_0 \circ \mathbf{b}_\alpha) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \} \times \{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_\alpha^T) \mathbf{h}_1^\psi + \mathbf{Z}(t) \mathbf{h}_1^\gamma \} d\Lambda_{s0}(t) \Big] d\mathbf{b}_\alpha \\
& + \sum_\alpha \int_{\mathbf{b}_\alpha} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \Delta_s h_2(V_s) \right. \\
& \quad \left. - \int_0^{V_s} h_2(t) \exp \{ \tilde{\mathbf{Z}}(t) (\boldsymbol{\psi} \circ \mathbf{b}_\alpha) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \} d\Lambda_{s0}(t) \right] d\mathbf{b}_\alpha, \tag{4.28}
\end{aligned}$$

where  $A'(D(t_j; \phi_0))$  and  $C'(Y_j; D(t_j; \phi_0))$  are the derivatives of  $A(D(t_j; \phi))$  and  $C(Y_j; D(t_j; \phi))$  with respect to  $\phi$  evaluated at  $\phi_0$  and  $B'(\boldsymbol{\beta}_0; \mathbf{b})$  is the derivative of  $B(\boldsymbol{\beta}; \mathbf{b})$  with respect to  $\boldsymbol{\beta}$  evaluated at  $\boldsymbol{\beta}_0$ . Note that (4.28) holds with probability one, so it may not hold for any  $V_s \in [0, \tau]$  when  $\Delta_s = 0$ . However, by the similar arguments done in Section 4.4.1, if we integrate both sides from  $V_s$  to  $\tau$  and subtract the obtained equation from (4.28) at  $\Delta_s = 0$  and  $V_s = \tau$ , it is easily shown that (4.28) also holds for any  $V_s \in [0, \tau]$  when  $\Delta_s = 0$ . Hence, the proof of the invertibility of  $\Omega$  will be completed if we can show  $\mathbf{h}_1 = 0$  and  $h_2(t) = 0$  from (4.28).

To show  $\mathbf{h}_1 = 0$ , particularly we let  $\Delta_s = 0$  and  $V_s = 0$  in (4.28) and obtain

$$\begin{aligned}
0 &= \sum_\alpha \int_{\mathbf{b}_\alpha} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \sum_{j=1}^{n_N} \frac{1}{A(D(t_j; \phi_0))} (Y_j \mathbf{X}_j - B'(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)) \mathbf{h}_1^\beta \right. \\
& \quad + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}_\alpha) - B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))^2} \right) A'(D(t_j; \phi_0)) + C'(Y_j; D(t_j; \phi_0)) \right\} \mathbf{h}_1^\phi \\
& \quad + \frac{1}{2} (\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0})^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} (\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0}) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \\
& \quad \left. + \left( \mathbf{b}_\alpha - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0} \right)^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_{\alpha 0}} \right] d\mathbf{b}_\alpha \\
&= \sum_\alpha \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right\} w_{\alpha 0} \right. \\
& \quad \times \left[ \mathbb{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \right] \times \left( \sum_{j=1}^{n_N} \frac{1}{A(D(t_j; \phi_0))} Y_j \mathbf{X}_j \mathbf{h}_1^\beta \right) \right.
\end{aligned}$$



$$\begin{aligned}
& - \frac{1}{A(D(t_j; \phi_0))} A'(D(t_j; \phi_0)) h_1^\phi \left( \mathbf{X}_j \boldsymbol{\beta}_0 \sum_{\alpha} \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right\} w_{\alpha 0} \right. \right. \\
& \quad \left. \left. \times E_{\mathbf{b}_{\alpha} | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \right] \right] \right] \\
& + \tilde{\mathbf{X}}_j \sum_{\alpha} \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right\} w_{\alpha 0} E_{\mathbf{b}_{\alpha} | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \mathbf{b}_{\alpha} \right] \right] \right] \Big\} \\
& = 0.
\end{aligned}$$

Based on assumption (A6),  $\mathbf{h}_1^\beta = 0$  and  $h_1^\phi = 0$ .

Then, we examine the constant terms without  $\mathbf{Y}$  in (4.29). Since  $\mathbf{h}_1^\beta = 0$  and  $h_1^\phi = 0$ , (4.29) becomes

$$\begin{aligned}
& \sum_{\alpha} \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right\} w_{\alpha 0} \right. \\
& \quad \times \left[ E_{\mathbf{b}_{\alpha} | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \right] \times \left( \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha 0} \right. \right. \\
& \quad \left. \left. - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_{\alpha}} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_{\alpha}} \right) \right. \\
& \quad \left. + E_{\mathbf{b}_{\alpha} | \alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \times \left( \frac{1}{2} \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_{\alpha} \right. \right. \right. \\
& \quad \left. \left. \left. - \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha} + \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_{\alpha}} \right) \right] \right] \Big\} \\
& = E_{\alpha, \mathbf{b}} \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} \right\} \times \exp \left\{ - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \right. \\
& \quad \times \left( \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha 0} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right. \\
& \quad \left. \left. - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_{\alpha}} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_{\alpha}} + \frac{1}{2} \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_{\alpha} - \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha} + \mathbf{b}_{\alpha}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_{\alpha}} \right) \right] \\
& = 0,
\end{aligned}$$

where  $\mathbf{b}$  follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_{b0}(\sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi_0)))) + \boldsymbol{\mu}_{\alpha 0}$  and covariance  $\boldsymbol{\Sigma}_{b0}$ . For any fixed  $\tilde{\mathbf{X}}$ , treating  $\mathbf{X}^T \mathbf{Y}$  as a parameter

in this normal family,  $\mathbf{b} = \sum_{k=1}^K I(\alpha = k) \mathbf{b}_k$  is the complete statistic for  $\mathbf{X}^T \mathbf{Y}$ , therefore,

$$\begin{aligned} & \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \times \left( \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha 0} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right. \\ & \quad \left. - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_\alpha} + \frac{1}{2} \mathbf{b}_\alpha^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_\alpha - \mathbf{b}_\alpha^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_\alpha + \mathbf{b}_\alpha^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} \right) \\ & = 0. \end{aligned}$$

Since  $\exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0} - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))} \right\} \neq 0$ , by (A6), we have

$$\begin{aligned} & \frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_{\alpha 0} - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) + \left( -\frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T + \mathbf{b}_\alpha^T \right) \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_\alpha} \\ & \quad + \frac{1}{2} \mathbf{b}_\alpha^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{b}_\alpha - \mathbf{b}_\alpha^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} \boldsymbol{\mu}_\alpha = 0. \end{aligned}$$

$\because -\frac{1}{2} \boldsymbol{\mu}_{\alpha 0}^T + \mathbf{b}_\alpha^T \neq 0$  and  $\boldsymbol{\Sigma}_{b0}^{-1} \neq 0 \implies \because \mathbf{h}_1^{\mu_\alpha} = 0, \alpha = 1, \dots, K$ , by (A6).

$\because 1/w_\alpha \neq 0 \implies \because \mathbf{h}_1^{w_\alpha} = 0, \alpha = 1, \dots, K$ , by (A6).

$\because \boldsymbol{\Sigma}_{b0}^{-1} \neq 0 \implies \because \mathbf{D}_b = \mathbf{0}$  by (A6).

Next, we let  $\Delta_s = 0$  in (4.28) and obtain

$$\begin{aligned} 0 &= \sum_{\alpha} \int_{\mathbf{b}_\alpha} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ \sum_{j=1}^{n_N} \frac{1}{A(D(t_j; \phi_0))} (Y_j \mathbf{X}_j - B'(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)) \mathbf{h}_1^\beta \right. \\ & \quad \left. + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \mathbf{b}_\alpha) - B(\boldsymbol{\beta}_0; \mathbf{b}_\alpha)}{A(D(t_j; \phi_0))^2} \right) A'(D(t_j; \phi_0)) + C'(Y_j; D(t_j; \phi_0)) \right\} \mathbf{h}_1^\phi \right. \\ & \quad \left. + \frac{1}{2} (\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0})^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_{b0}^{-1} (\mathbf{b}_\alpha - \boldsymbol{\mu}_{\alpha 0}) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_{b0}^{-1} \mathbf{D}_b) \right. \\ & \quad \left. + \left( \mathbf{b}_\alpha - \frac{1}{2} \boldsymbol{\mu}_{\alpha 0} \right)^T \boldsymbol{\Sigma}_{b0}^{-1} \mathbf{h}_1^{\mu_\alpha} + \frac{1}{w_{\alpha 0}} \mathbf{h}_1^{w_{\alpha 0}} \right. \\ & \quad \left. - \int_0^{V_s} \exp \left\{ \tilde{\mathbf{Z}}(t) (\boldsymbol{\psi}_0 \circ \mathbf{b}_\alpha) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \right\} \times \left\{ (\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_\alpha^T) \mathbf{h}_1^\psi + \mathbf{Z}(t) \mathbf{h}_1^\gamma \right\} d\Lambda_{s0}(t) \right] d\mathbf{b}_\alpha \\ & \quad + \sum_{\alpha} \int_{\mathbf{b}_\alpha} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \times \left[ - \int_0^{V_s} h_2(t) \exp \left\{ \tilde{\mathbf{Z}}(t) (\boldsymbol{\psi} \circ \mathbf{b}_\alpha) + \mathbf{Z}(t) \boldsymbol{\gamma}_0 \right\} d\Lambda_{s0}(t) \right] d\mathbf{b}_\alpha. \end{aligned}$$

Since  $\mathbf{h}_1^\beta = 0$ ,  $\mathbf{h}_1^\phi$ ,  $\mathbf{h}_1^{\mu_\alpha} = 0$ ,  $\mathbf{h}_1^{w_{\alpha 0}} = 0$ ,  $\alpha = 1, \dots, K$ , and  $\mathbf{D}_b = 0$ , the above expression can

be written as

$$0 = E_{\alpha, \mathbf{b}} \left[ \exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0}) - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \right. \\ \left. \times \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}_{\alpha}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} \times [(\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_{\alpha}^T) \mathbf{h}_1^{\psi} + \mathbf{Z}(t) \mathbf{h}_1^{\gamma} + h_2(t)] d\Lambda_{s0}(t) \right], \quad (4.30)$$

where  $\mathbf{b}_{\alpha}$  follows a multivariate normal distribution with mean  $\boldsymbol{\Sigma}_{b0} [\sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{Z}}_j / A(D(t_j; \phi_0)))] + \boldsymbol{\mu}_{\alpha 0}$  and covariance  $\boldsymbol{\Sigma}_{b0}$ . Likewise, for any fixed  $\tilde{\mathbf{X}}$ , treating  $\mathbf{X}^T \mathbf{Y}$  as a parameter in this normal family,  $\mathbf{b}_{\alpha}$  is the complete statistic for  $\mathbf{X}^T \mathbf{Y}$ , therefore,

$$\exp \left\{ \sum_{j=1}^{n_N} \frac{Y_j}{A(D(t_j; \phi_0))} (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0}) - \sum_{j=1}^{n_N} \frac{B(\boldsymbol{\beta}_0; \mathbf{b}_{\alpha})}{A(D(t_j; \phi_0))} \right\} \\ \times \int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}_{\alpha}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} \times [(\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_{\alpha}^T) \mathbf{h}_1^{\psi} + \mathbf{Z}(t) \mathbf{h}_1^{\gamma} + h_2(t)] d\Lambda_{s0}(t) \\ = 0.$$

Since  $\exp \{ \sum_{j=1}^{n_N} [Y_j (\mathbf{X}_j \boldsymbol{\beta}_0 + \tilde{\mathbf{X}}_j \boldsymbol{\mu}_{\alpha 0}) / A(D(t_j; \phi_0))] - \sum_{j=1}^{n_N} [B(\boldsymbol{\beta}_0; \mathbf{b}) / A(D(t_j; \phi_0))] \} \neq 0$ , equivalently

$$\int_0^{V_s} \exp \{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi}_0 \circ \mathbf{b}_{\alpha}) + \mathbf{Z}(t)\boldsymbol{\gamma}_0 \} \times [(\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_{\alpha}^T) \mathbf{h}_1^{\psi} + \mathbf{Z}(t) \mathbf{h}_1^{\gamma} + h_2(t)] d\Lambda_{s0}(t) = 0$$

by assumption (A6). From assumption (A6), this immediately gives  $\mathbf{h}_1^{\psi} = 0$ ,  $\mathbf{h}_1^{\gamma} = 0$  and  $h_2(t) = 0$ . Hence, the proof of condition (c) is completed.

Since the conditions (a)–(d) have been proved, Theorem 3.3.1 of van der Vaart and Wellner (1996) concludes that  $\sqrt{m_s}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})$  weakly converges to a tight random element in  $\ell^{\infty}(\mathcal{H})$ . Furthermore, we obtain



$$\begin{aligned}
& \sqrt{m_s} \nabla U_{\xi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\
&= \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P})\{l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2]\} + o_P(1),
\end{aligned} \tag{4.31}$$

where  $o_P(1)$  is a random variable which converges to zero in probability in  $\ell^\infty(\mathcal{H})$ .

On the other hand, from (4.27), we have

$$\begin{aligned}
& \sqrt{m_s} \nabla U_{\xi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\
&= \sqrt{m_s}\{(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\widehat{\Lambda}_s - \Lambda_{s0})(t)\}.
\end{aligned} \tag{4.32}$$

By denoting  $(\mathbf{h}_1^*, h_2^*) = \Omega^{-1}(\mathbf{h}_1, h_2)$ , we have  $(\mathbf{h}_1, h_2) = \Omega(\mathbf{h}_1^*, h_2^*)$ , and replacing  $(\mathbf{h}_1, h_2)$  with  $(\mathbf{h}_1^*, h_2^*)$  in (4.31) and (4.32) leads to the followings, respectively.

$$\begin{aligned}
& \sqrt{m_s} \nabla U_{\xi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1^*, h_2^*] \\
&= \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P})\{l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]\} + o_P(1),
\end{aligned}$$

and

$$\begin{aligned}
& \sqrt{m_s} \nabla U_{\xi_0}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})[\mathbf{h}_1^*, h_2^*] \\
&= \sqrt{m_s}\left\{(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1^*, h_2^*] + \int_0^\tau \Omega_2[\mathbf{h}_1^*, h_2^*] d(\widehat{\Lambda}_s - \Lambda_{s0})(t)\right\} \\
&= \sqrt{m_s}\left\{(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\widehat{\Lambda}_s - \Lambda_{s0})(t)\right\}.
\end{aligned}$$

Thus, we obtain

$$\begin{aligned}
& \sqrt{m_s}\left\{(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^T \mathbf{h}_1 + \int_0^\tau h_2(t) d(\widehat{\Lambda}_s - \Lambda_{s0})(t)\right\} \\
&= \sqrt{m_s}(\mathbf{P}_{m_s} - \mathbf{P})\{l_{\theta}(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]\} + o_P(1).
\end{aligned} \tag{4.33}$$

Note that the first term on the right-hand side in (4.33) is  $\sqrt{m_s}\{U_{m_s}(\boldsymbol{\theta}_0, \Lambda_{s0}) - U_s(\boldsymbol{\theta}_0, \Lambda_{s0})\}$ , which is an empirical process in the space  $\ell^\infty(\mathcal{H})$ , and it is shown that  $\mathcal{G}$  is P-Donsker in Section 4.4.3.2. Therefore,  $\sqrt{m_s}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0, \widehat{\Lambda}_s - \Lambda_{s0})$  weakly converges to a Gaussian process in  $\ell^\infty(\mathcal{H})$ .

In particular, if we choose  $h_2 = 0$  in (4.33), then  $\widehat{\boldsymbol{\theta}}^T \mathbf{h}_1$  is an asymptotic linear estimator for  $\boldsymbol{\theta}_0^T \mathbf{h}_1$  with influence function being  $l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1^* + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2^*]$ . Since this influence function is in the linear space spanned by the score functions for  $\boldsymbol{\theta}_0$  and  $\Lambda_{s0}$ , Proposition 3.3.1 (p65) in Bickel, Klaassen, Ritov and Wellner (1993) concludes that the influence function is the same as the efficient influence function for  $\boldsymbol{\theta}_0^T \mathbf{h}_1$ ; that is  $\widehat{\boldsymbol{\theta}}$  is an efficient estimator for  $\boldsymbol{\theta}_0$  and Theorem 4.2 is proved.

### 4.4.3 Supplementary proofs

The proofs for P-Donsker property of the classes  $\mathcal{F}$  and  $\mathcal{G}$  needed in Sections 4.4.1 and 4.4.2 are presented in Sections 4.4.3.1~4.4.3.2 respectively. In Section 4.4.3.3, we prove *Fréchet* differentiability of  $U(\xi)$  at  $\xi_0$  and derive the derivative operator  $\nabla U_{\xi_0}$  used in Section 4.4.2.

#### 4.4.3.1 Proof of P-Donsker property of $\mathcal{F}$

We defined that a class  $\mathcal{F} = \{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : v \in [0, \tau], \boldsymbol{\theta} \in \Theta, \Lambda_s \in \mathcal{A}, s = 1, \dots, S\}$ , where  $\mathcal{A} = \{\Lambda_s \in \mathbb{W}, \Lambda_s(0) = 0, \Lambda_s(\tau) \leq B_{s0}, s = 1, \dots, S\}$ ,  $B_{s0}$  is the constant given in the second step and  $\mathbb{W}$  contains all nondecreasing functions in  $[0, \tau]$ . We can rewrite  $Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)$  as

$$Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) = Q_1(v, \mathbf{O}; \boldsymbol{\theta}) \frac{Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)}{Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)},$$

where

$$\begin{aligned}
& Q_1(v, \mathbf{O}; \boldsymbol{\theta}) \\
&= \exp \left\{ \mathbf{Z}(v) \boldsymbol{\gamma} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\Sigma}_b (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T + \frac{1}{2} R(v) \right\}, \\
& Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta+1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right. \\
&\quad \left. - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] + R(t) \right\} d\Lambda_s(t) \right\} w_{\alpha} d\mathbf{b}_{\alpha}, \\
& Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right. \\
&\quad \left. - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right\} d\Lambda_s(t) \right\} w_{\alpha} d\mathbf{b}_{\alpha},
\end{aligned}$$

$R(t) = (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T)^T$ ,  $R(v)$  is  $R(t)$  evaluated at  $t = v$ ,  $B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha}) = B(\boldsymbol{\beta}; g_1(\mathbf{b}_{\alpha}))$ ,  $B_2(\boldsymbol{\beta}; \mathbf{b}_{\alpha}) = B(\boldsymbol{\beta}; g_2(\mathbf{b}_{\alpha}))$ ,  $g_1(\mathbf{b}_{\alpha}) = \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \boldsymbol{\Sigma}_b \left[ \sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi))) + (\Delta+1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right]^T + \boldsymbol{\mu}_{\alpha}$  and  $g_2(\mathbf{b}_{\alpha}) = \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \boldsymbol{\Sigma}_b \left[ \sum_{j=1}^{n_N} (Y_j \tilde{\mathbf{X}}_j / A(D(t_j; \phi))) + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right]^T + \boldsymbol{\mu}_{\alpha}$ .

Using assumption (A3), we can easily show that  $Q_1(v, \mathbf{O}; \boldsymbol{\theta})$  is continuously differentiable with respect to  $v$  and  $\boldsymbol{\theta}$ , and

$$\|\nabla_{\boldsymbol{\theta}} Q_1(v, \mathbf{O}; \boldsymbol{\theta})\| + \left| \frac{d}{dv} Q_1(v, \mathbf{O}; \boldsymbol{\theta}) \right| \leq e^{k_1 + k_2 \|\mathbf{Y}\|}$$

for some positive constants  $k_1$  and  $k_2$ . Furthermore, it holds that

$$\begin{aligned}
& \|\nabla_{\boldsymbol{\theta}} Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \\
& \leq \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \left[ \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
& \quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \times e^{k_3 \|\mathbf{b}_{\alpha}\| + k_4 \|\mathbf{Y}\| + k_5(\alpha)} \times B_{s_0} \times w_{\alpha} \right] d\mathbf{b}_{\alpha} \\
& \leq e^{k_6 + k_7 \|\mathbf{Y}\|}
\end{aligned}$$

$$\text{and} \quad \|\nabla_{\boldsymbol{\theta}} Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \leq e^{k_8 + k_9 \|\mathbf{Y}\|}$$

for some positive constants  $k_3, k_4, k_6, k_7, k_8$ , and  $k_9$ , and a deterministic function of  $\alpha$ ,  $k(\alpha)$ . Additionally, note that, for any  $0 < \Lambda < \infty$ ,  $0 < e^{-\Lambda} < 1$  and  $e^{-\Lambda} < \Lambda$  and thus  $e^{-\Lambda_1} - e^{-\Lambda_2} < \Lambda_1 - \Lambda_2$  for any  $\Lambda_1$  and  $\Lambda_2$  over  $(0, \infty)$ . Hence,

$$\begin{aligned}
& |Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s1}) - Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s2})| \\
& = \left| \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
& \quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \left[ \exp \left\{ -\int_0^v \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \right. \\
& \quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right. \right. \right. \\
& \quad \left. \left. + R(t) \right\} d\Lambda_{s1}(t) \right] - \exp \left\{ -\int_0^v \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
& \quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right. \right. \\
& \quad \left. \left. + R(t) \right\} d\Lambda_{s2}(t) \right] w_{\alpha} d\mathbf{b}_{\alpha} \right| \\
& \leq \left| \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
& \quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \left[ \int_0^v \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \right.
\end{aligned}$$

$$\begin{aligned}
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_\alpha \right] \\
& + R(t) \Big\} d(\Lambda_{s1} - \Lambda_{s2})(t) \Big] w_\alpha d\mathbf{b}_\alpha \Big| \\
= & \left| \sum_\alpha \int_{\mathbf{b}_\alpha} \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta+1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_\alpha - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} \right\} \times (2\pi)^{d_b/2} \right. \\
& \times (2\pi)^{-d_b/2} \times \left[ \int_0^v \exp \left\{ -\frac{1}{2} [\mathbf{b}_\alpha - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T]^T [\mathbf{b}_\alpha - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T] \right\} \right. \\
& \times \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \\
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_\alpha \right] \\
& \left. \left. \left. + R(t) \right\} d(\Lambda_{s1} - \Lambda_{s2})(t) \right] w_\alpha d\mathbf{b}_\alpha \right| \\
\leq & \left| \sum_\alpha \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta+1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_\alpha \right\} \times (2\pi)^{d_b/2} \right. \right. \\
& \times \int_0^v \left[ \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_\alpha \right] + R(t) \Big\} \\
& \times (2\pi)^{-d_b/2} \int_{\mathbf{b}_\alpha} \exp \left\{ -\sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} \right\} \\
& \times \exp \left\{ -\frac{1}{2} [\mathbf{b}_\alpha - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T]^T [\mathbf{b}_\alpha - ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T] \right\} d\mathbf{b}_\alpha \Big] \\
& \left. d(\Lambda_{s1} - \Lambda_{s2})(t) \right] w_\alpha \Big| \\
= & \left| \sum_\alpha \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta+1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_\alpha \right\} \times (2\pi)^{d_b/2} \right. \right. \\
& \times \int_0^v \left[ \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_\alpha \right] + R(t) \Big\} \\
& \times \mathbb{E}_{\mathbf{b}_\alpha | \alpha} \left[ \exp \left\{ -\sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} \right\} \right] \Big] d(\Lambda_{s1} - \Lambda_{s2})(t) \Big] w_\alpha \Big|
\end{aligned}$$

$$\begin{aligned}
&= \left| \sum_{\alpha} \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right\} \times (2\pi)^{d_b/2} \right. \right. \\
&\quad \times \left[ (\Lambda_{s1}(v) - \Lambda_{s2}(v)) \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] + R(v) \right\} \right. \\
&\quad \left. \left. \times \mathbb{E}_{\mathbf{b}_{\alpha}|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \right] \right. \right. \\
&\quad \left. \left. - \int_0^v \left[ (\Lambda_{s1}(t) - \Lambda_{s2}(t)) \frac{d}{dt} \left[ \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \right. \right. \\
&\quad \left. \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] + R(t) \right\} \right. \right. \right. \\
&\quad \left. \left. \left. \times \mathbb{E}_{\mathbf{b}_{\alpha}|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \right] \right] \right] dt \right] w_{\alpha} \right| \\
&\leq \sum_{\alpha} w_{\alpha} \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right\} \times (2\pi)^{d_b/2} \right. \\
&\quad \times \left[ |\Lambda_{s1}(v) - \Lambda_{s2}(v)| \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] + R(v) \right\} \right. \\
&\quad \left. \left. \times \mathbb{E}_{\mathbf{b}_{\alpha}|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \right] \right. \right. \\
&\quad \left. \left. + \int_0^v \left[ |\Lambda_{s1}(t) - \Lambda_{s2}(t)| \left| \frac{d}{dt} \left[ \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \right. \right. \right. \\
&\quad \left. \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] + R(t) \right\} \right. \right. \right. \\
&\quad \left. \left. \left. \times \mathbb{E}_{\mathbf{b}_{\alpha}|\alpha} \left[ \exp \left\{ - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \right] \right] \right] dt \right] \\
&= |\Lambda_{s1}(v) - \Lambda_{s2}(v)| \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2}) ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \\
&\quad \left. + (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right)^T + R(v) \right\} \times (2\pi)^{d_b/2} \\
&\quad \times \mathbb{E}_{\alpha, \mathbf{b}} \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \right]
\end{aligned}$$

$$\begin{aligned}
& + \int_0^v \left[ |\Lambda_{s1}(t) - \Lambda_{s2}(t)| \times (2\pi)^{d_b/2} \right. \\
& \times \mathbb{E}_{\alpha, \mathbf{b}} \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_\alpha \right\} \right. \\
& \times \left| \frac{d}{dt} \left[ \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \Sigma_b^{1/2}) ((\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \Sigma_b^{1/2})^T + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \\
& + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \Sigma_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_\alpha \right] \right. \\
& \left. \left. \left. + R(t) - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} \right\} \right] \right] dt \\
& \leq (2\pi)^{d_b/2} \exp \left\{ \frac{1}{2} ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \Sigma_b^{1/2}) ((\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \Sigma_b^{1/2})^T + \mathbf{Z}(v) \boldsymbol{\gamma} \right. \\
& \left. + (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \Sigma_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right)^T + R(v) \right\} \\
& \times \mathbb{E}_{\alpha, \mathbf{b}} \left[ \exp \left\{ \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_\alpha - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_\alpha)}{A(D(t_j; \phi))} \right\} \right] \\
& \times \left[ |\Lambda_{s1}(v) - \Lambda_{s2}(v)| + \int_0^v |\Lambda_{s1}(t) - \Lambda_{s2}(t)| dt \right] \\
& \leq e^{k_{10} + k_{11} \|\mathbf{Y}\|} \left[ |\Lambda_{s1}(v) - \Lambda_{s2}(v)| + \int_0^\tau |\Lambda_{s1}(t) - \Lambda_{s2}(t)| dt \right],
\end{aligned}$$

where  $k_{10}$  and  $k_{11}$  are positive constants. Similarly,

$$\begin{aligned}
& |Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s1}) - Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_{s2})| \\
& \leq e^{k_{12} + k_{13} \|\mathbf{Y}\|} \left[ |\Lambda_{s1}(v) - \Lambda_{s2}(v)| + \int_0^\tau |\Lambda_{s1}(t) - \Lambda_{s2}(t)| dt \right],
\end{aligned}$$

where  $k_{12}$  and  $k_{13}$  are positive constants.

On the other hand, there exist positive constants  $k_{14}, \dots, k_{26}$  such that

$$\begin{aligned}
& |Q_1(v, \mathbf{O}; \boldsymbol{\theta})| \\
& = \left| \exp \left\{ \mathbf{Z}(v) \boldsymbol{\gamma} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \Sigma_b (\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T)^T + \frac{1}{2} R(v) \right\} \right|
\end{aligned}$$

$$\leq e^{k_{14}+k_{15}\|\mathbf{Y}\|},$$

$$\begin{aligned}
& |Q_2(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)| \\
&= \left| \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
&\quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right. \right. \\
&\quad \left. \left. + R(t) \right\} d\Lambda_s(t) \right\} w_{\alpha} d\mathbf{b}_{\alpha} \Big| \\
&\leq \left| \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
&\quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \times \left[ 2 \int_0^v \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right. \right. \\
&\quad \left. \left. + R(t) \right\} d\Lambda_s(t) \right] w_{\alpha} d\mathbf{b}_{\alpha} \Big| \\
&\leq \left| \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + (\Delta + 1)(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} \right. \right. \\
&\quad \left. \left. - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right\} \times 2 \times \exp \{ k_{16} \|\mathbf{b}_{\alpha}\| + k_{17} \|\mathbf{Y}\| + k_{18} \|\boldsymbol{\mu}_{\alpha}\| + k_{19} \} \times B_{s0} \times w_{\alpha} d\mathbf{b}_{\alpha} \right| \\
&\leq e^{k_{19}+k_{20}\|\mathbf{Y}\|},
\end{aligned}$$

and

$$\begin{aligned}
& Q_3(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right. \\
&\quad \left. - \int_0^{V_s} \exp \left\{ (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \boldsymbol{\Sigma}_b^{1/2} \mathbf{b}_{\alpha} + \mathbf{Z}(t) \boldsymbol{\gamma} \right. \right. \\
&\quad \left. \left. + (\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \left[ \boldsymbol{\Sigma}_b \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(t) \circ \boldsymbol{\psi}^T) \right)^T + \boldsymbol{\mu}_{\alpha} \right] \right\} d\Lambda_s(t) \right\} w_{\alpha} d\mathbf{b}_{\alpha}
\end{aligned}$$



$$\begin{aligned}
&\geq \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} \exp \left\{ -\frac{1}{2} \mathbf{b}_{\alpha}^T \mathbf{b}_{\alpha} + \left( \sum_{j=1}^{n_N} \frac{Y_j \tilde{\mathbf{X}}_j}{A(D(t_j; \phi))} + \Delta(\tilde{\mathbf{Z}}(v) \circ \boldsymbol{\psi}^T) \right) \boldsymbol{\mu}_{\alpha} - \sum_{j=1}^{n_N} \frac{B_1(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))} \right. \\
&\quad \left. - \exp \{ k_{22} \|\mathbf{b}_{\alpha}\| + k_{23} \|\mathbf{Y}\| + k_{24} \|\boldsymbol{\mu}_{\alpha}\| + k_{25} \} \times B_{s0} \right\} w_{\alpha} d\mathbf{b}_{\alpha} \\
&\geq k_{26} > 0.
\end{aligned}$$

Moreover,

$$\begin{aligned}
&\|\nabla_{\boldsymbol{\theta}} Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)\| + \left| \frac{d}{dv} Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \right| \\
&= \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left( \nabla_{\boldsymbol{\theta}} \frac{Q_2}{Q_3} \right) \right\| + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left( \frac{d}{dv} \left( \frac{Q_2}{Q_3} \right) \right) \right\| \\
&= \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left[ \left( \nabla_{\boldsymbol{\theta}} Q_2 \right) \frac{1}{Q_3} + Q_2 \frac{(-1)}{Q_3^2} (\nabla_{\boldsymbol{\theta}} Q_3) \right] \right\| \\
&\quad + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + Q_1 \left[ \left( \frac{d}{dv} Q_2 \right) \frac{1}{Q_3} + Q_2 \frac{(-1)}{Q_3^2} \left( \frac{d}{dv} Q_3 \right) \right] \right\| \\
&= \left\| \left( \nabla_{\boldsymbol{\theta}} Q_1 \right) \frac{Q_2}{Q_3} + \left( \nabla_{\boldsymbol{\theta}} Q_2 \right) \frac{Q_1}{Q_3} - \left( \nabla_{\boldsymbol{\theta}} Q_3 \right) \frac{Q_1 Q_2}{Q_3^2} \right\| \\
&\quad + \left\| \left( \frac{d}{dv} Q_1 \right) \frac{Q_2}{Q_3} + \left( \frac{d}{dv} Q_2 \right) \frac{Q_1}{Q_3} - \left( \frac{d}{dv} Q_3 \right) \frac{Q_1 Q_2}{Q_3^2} \right\| \\
&\leq \left( \|\nabla_{\boldsymbol{\theta}} Q_1\| + \left| \frac{d}{dv} Q_1 \right| \right) \left| \frac{Q_2}{Q_3} \right| + \left( \|\nabla_{\boldsymbol{\theta}} Q_2\| + \left| \frac{d}{dv} Q_2 \right| \right) \left| \frac{Q_1}{Q_3} \right| + \left( \|\nabla_{\boldsymbol{\theta}} Q_3\| + \left| \frac{d}{dv} Q_3 \right| \right) \left| \frac{Q_1 Q_2}{Q_3^2} \right| \\
&\leq e^{k_{27} + k_{28} \|\mathbf{Y}\|},
\end{aligned}$$

for some positive constants  $k_{27}$  and  $k_{28}$ . Therefore, by the mean-value theorem, we conclude that, for any  $(v_1, \boldsymbol{\theta}_1, \Lambda_{s1})$  and  $(v_2, \boldsymbol{\theta}_2, \Lambda_{s2})$  in  $[0, \tau] \times \Theta \times \mathcal{A}$ ,

$$\begin{aligned}
&|Q(v_1, \mathbf{O}; \boldsymbol{\theta}_1, \Lambda_{s1}) - Q(v_2, \mathbf{O}; \boldsymbol{\theta}_2, \Lambda_{s2})| \\
&\leq e^{k_{27} + k_{28} \|\mathbf{Y}\|} \left[ \|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2\| + |\Lambda_{s1}(V) - \Lambda_{s2}(V)| + \int_0^{\tau} |\Lambda_{s1}(t) - \Lambda_{s2}(t)| dt + |v_1 - v_2| \right]
\end{aligned} \tag{4.34}$$

holds for some positive constants  $k_{27}$  and  $k_{28}$  and  $0 \leq V \leq \tau$  ( $V = v_1$  or  $v_2$ ).

Applying Theorem 2.7.5 (p159) in van der Vaart and Wellner (1996) to our situation, the entropy number for the class  $\mathcal{A}$  satisfies  $\log N_{[\cdot]}(\varepsilon, \mathcal{A}, L_2(P)) \leq K/\varepsilon$ , where  $K$  is a constant. Thus, we can find  $\exp\{K/\varepsilon\}$  brackets,  $\{[L_j, U_j]\}$ , to cover the class  $\mathcal{A}$  such that  $\|U_j - L_j\|_{L_2(P)} \leq \varepsilon$  for each pair of  $[L_j, U_j]$ . On the other hand, we can further find a partition of  $[0, \tau] \times \Theta$ , say  $I_1 \cup I_2 \cup \dots$ , such that the number of partitions is of the order  $(1/\varepsilon)^{d_\theta+1}$ , and, for any  $(v_1, \boldsymbol{\theta}_1)$  and  $(v_2, \boldsymbol{\theta}_2)$  in the same partition, their Euclidean distance is less than  $\varepsilon$ . Therefore, the partition  $\{I_1, I_2, \dots\} \times \{[L_j, U_j]\}$  bracket covers  $[0, \tau] \times \Theta \times \mathcal{A}$ , and the total number of the partitions is of order  $(1/\varepsilon)^{d_\theta+1} \exp\{1/\varepsilon\}$ . Hence, from (4.34), for any  $I_l$  and  $[L_j, U_j]$ , the set of the functions  $\{Q(v, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) : (v, \boldsymbol{\theta}) \in I_l, \Lambda_s \in \mathcal{A}, \Lambda_s \in [L_j, U_j]\}$  can be bracket covered by

$$\left[ Q(v_l, \mathbf{O}; \boldsymbol{\theta}_l, \Lambda_{sl}) - e^{k_{27}+k_{28}\|\mathbf{Y}\|} \left\{ \varepsilon + |U_j(V) - L_j(V)| + \int_0^\tau |U_j(t) - L_j(t)| dt \right\}, \right. \\ \left. Q(v_l, \mathbf{O}; \boldsymbol{\theta}_l, \Lambda_{sl}) + e^{k_{27}+k_{28}\|\mathbf{Y}\|} \left\{ \varepsilon + |U_j(V) - L_j(V)| + \int_0^\tau |U_j(t) - L_j(t)| dt \right\} \right], \quad (4.35)$$

where  $(v_l, \boldsymbol{\theta}_l)$  is a fixed point in  $I_l$  and  $\Lambda_{sl}$  is a fixed function in  $[L_j, U_j]$ . Note that the  $L_2(P)$  distance between these two functions in the above bracket (4.35) is less than  $O(\varepsilon)$ . Therefore, we have

$$N_{[\cdot]}(\varepsilon, \mathcal{F}, \|\cdot\|_{L_2(P)}) \leq O\left(\left(\frac{1}{\varepsilon}\right)^{d_\theta+1} e^{1/\varepsilon}\right).$$

Furthermore,  $\mathcal{F}$  has an  $L_2(P)$ -integrable covering function, which is equal to  $O(e^{k_{27}+k_{28}\|\mathbf{Y}\|})$ . From Theorem 2.5.6 (p130) in van der Vaart and Wellner (1996),  $\mathcal{F}$  is P-Donsker.

Additionally, in the above derivation, we also note that all the functions in  $\mathcal{F}$  are bounded from below by  $e^{-k_{29}-k_{30}\|\mathbf{Y}\|}$  for some positive constants  $k_{29}$  and  $k_{30}$ .

#### 4.4.3.2 Proof of P-Donsker property of $\mathcal{G}$

Recall that we defined the class

$$\begin{aligned}\mathcal{G} &= \left\{ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2], \right. \\ &\quad \left. \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \leq \delta, (\mathbf{h}_1, h_2) \in \mathcal{H} \right\},\end{aligned}$$

where  $(\mathbf{h}_1^\beta, \mathbf{h}_1^\phi, \mathbf{h}_1^{\Sigma_b}, \mathbf{h}_1^\mu, \mathbf{h}_1^w, \mathbf{h}_1^\psi, \mathbf{h}_1^\gamma)$  denote the corresponding components of  $\mathbf{h}_1$  for the parameters  $(\boldsymbol{\beta}, \boldsymbol{\phi}, \text{Vec}(\boldsymbol{\Sigma}_b), \boldsymbol{\mu}, \mathbf{w}, \boldsymbol{\psi}, \boldsymbol{\gamma})$ , respectively. We can write that for  $(\mathbf{h}_1, h_2) \in \mathcal{H}$ ,

$$\begin{aligned}& l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] \\ &= \left[ \rho_1(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 - \int_0^{V_s} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 d\Lambda_s(t) \right] + \Delta h_2(V_s) - \int_0^{V_s} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) h_2(t) d\Lambda_s(t),\end{aligned}$$

where

$$\begin{aligned}& \rho_1(\mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 \\ &= \left\{ \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b}_{\alpha} \right\}^{-1} \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\ &\quad \times \left[ \sum_{j=1}^{n_N} \frac{1}{A(D(t_j; \phi))} (Y_j \mathbf{X}_j - B'(\boldsymbol{\beta}; \mathbf{b}_{\alpha})) \mathbf{h}_1^{\beta} \right. \\ &\quad + \sum_{j=1}^{n_N} \left\{ - \left( \frac{Y_j(\mathbf{X}_j \boldsymbol{\beta} + \tilde{\mathbf{X}}_j \mathbf{b}_{\alpha}) - B(\boldsymbol{\beta}; \mathbf{b}_{\alpha})}{A(D(t_j; \phi))^2} \right) A'(D(t_j; \phi)) + C'(Y_j; D(t_j; \phi)) \right\} \mathbf{h}_1^{\phi} \\ &\quad + \frac{1}{2} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha})^T \boldsymbol{\Sigma}_b^{-1} \mathbf{D}_b \boldsymbol{\Sigma}_b^{-1} (\mathbf{b}_{\alpha} - \boldsymbol{\mu}_{\alpha}) - \frac{1}{2} \text{Tr}(\boldsymbol{\Sigma}_b^{-1} \mathbf{D}_b) \\ &\quad \left. + \left( \mathbf{b}_{\alpha} - \frac{1}{2} \boldsymbol{\mu}_{\alpha} \right)^T \boldsymbol{\Sigma}_b^{-1} \mathbf{h}_1^{\mu_{\alpha}} + \frac{1}{w_{\alpha}} \mathbf{h}_1^{w_{\alpha}} + \Delta_s \{ (\tilde{\mathbf{Z}}(V_s) \circ \mathbf{b}_{\alpha}^T) \mathbf{h}_1^{\psi} + \mathbf{Z}(V_s) \mathbf{h}_1^{\gamma} \} \right] d\mathbf{b}_{\alpha}, \\ & \rho_2(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 \\ &= \left\{ \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b}_{\alpha} \right\}^{-1} \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\ &\quad \times \exp \left\{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}_{\alpha}) + \mathbf{Z}(t) \boldsymbol{\gamma} \right\} \times \left[ (\tilde{\mathbf{Z}}(t) \circ \mathbf{b}_{\alpha}^T) \mathbf{h}_1^{\psi} + \mathbf{Z}(t) \mathbf{h}_1^{\gamma} \right] d\mathbf{b}_{\alpha},\end{aligned}$$

$$\begin{aligned}
& \rho_3(t, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \\
&= \left\{ \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\mathbf{b}_{\alpha} \right\}^{-1} \\
& \quad \times \sum_{\alpha} \int_{\mathbf{b}_{\alpha}} G(\mathbf{b}, \alpha, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) \times \exp \left\{ \tilde{\mathbf{Z}}(t)(\boldsymbol{\psi} \circ \mathbf{b}_{\alpha}) + \mathbf{Z}(t)\boldsymbol{\gamma} \right\} d\mathbf{b}_{\alpha},
\end{aligned}$$

$B'(\boldsymbol{\beta}; \mathbf{b})$  is the derivative of  $B(\boldsymbol{\beta}; \mathbf{b})$  with respect to  $\boldsymbol{\beta}$ ,  $A'(D(t_j; \phi))$  and  $C'(Y_j; D(t_j; \phi))$  are the derivatives of  $A(D(t_j; \phi))$  and  $C(Y_j; D(t_j; \phi))$  with respect to  $\phi$  respectively, and  $\mathbf{D}_b$  is the symmetric matrix such that  $\text{Vec}(\mathbf{D}_b) = \mathbf{h}_1^b$ .

For  $l = 1, 2, 3$ , we denote  $\nabla_{\boldsymbol{\theta}} \rho_l$  and  $\nabla_{\Lambda_s} \rho_l[\delta \Lambda_s]$  as the derivatives of  $\rho_l$  with respect to  $\boldsymbol{\theta}$  and  $\Lambda_s$  along the path  $\Lambda_s + \varepsilon \delta \Lambda_s$ . Then, using the similar arguments done in Section 4.4.3.1, it is verified that  $\nabla_{\Lambda_s} \rho_l[\delta \Lambda_s] = \int_0^t \rho_{l+3}(u, \mathbf{O}; \boldsymbol{\theta}, \Lambda_s) d\delta \Lambda_s(u)$  and there exist two positive constants  $q_1$  and  $q_2$  such that

$$\sum_l \{|\rho_l| + |\nabla_{\boldsymbol{\theta}} \rho_l|\} \leq e^{q_1 + q_2 \|\mathbf{Y}\|}$$

By the mean value theorem, we have that, for any  $(\boldsymbol{\theta}, \Lambda_s, \mathbf{h}_1, h_2)$  and  $(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s, \tilde{\mathbf{h}}_1, \tilde{h}_2)$  in  $\Xi \times \mathcal{H}$ ,

$$\begin{aligned}
& l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2] \\
&= l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \mathbf{h}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[h_2] \\
& \quad + l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2] \\
&= [l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T] \mathbf{h}_1 + [l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)][h_2] \\
& \quad + l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)([h_2] - [\tilde{h}_2]) \\
&= (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \left[ \frac{d}{d\boldsymbol{\theta}} l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right] \mathbf{h}_1 + \left[ \frac{d}{d\Lambda_s} l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right]^T [\Lambda_s - \tilde{\Lambda}_s] \mathbf{h}_1 \\
& \quad + (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \left[ \frac{d}{d\boldsymbol{\theta}} l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right] [h_2] + \left[ \frac{d}{d\Lambda_s} l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}^*, \Lambda_s=\Lambda_s^*} \right]^T [\Lambda_s - \tilde{\Lambda}_s][h_2] \\
& \quad + l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) + l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)([h_2] - [\tilde{h}_2])
\end{aligned}$$

$$\begin{aligned}
&= (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \nabla_{\boldsymbol{\theta}} \rho_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \mathbf{h}_1 - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d\Lambda_s^*(t) \mathbf{h}_1 \\
&\quad + \int_0^{V_s} \rho_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \tilde{\Lambda}_s)(t) \\
&\quad + \int_0^{V_s} \int_0^t \rho_5(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \tilde{\Lambda}_s)(u) \mathbf{h}_1 d\Lambda_s^*(t) \\
&\quad - \int_0^{V_s} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T (\Lambda_s - \tilde{\Lambda}_s) \mathbf{h}_1 dt - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
&\quad + \int_0^{V_s} \int_0^t \rho_6(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) d(\Lambda_s - \tilde{\Lambda}_s)(u) h_2(t) d\Lambda_s^*(t) \\
&\quad - \int_0^{V_s} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T (\Lambda_s - \tilde{\Lambda}_s)(t) h_2(t) dt \\
&\quad + \rho_1(\mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) - \int_0^{V_s} \rho_2(t, \mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T (\mathbf{h}_1 - \tilde{\mathbf{h}}_1) d\tilde{\Lambda}_s(t) \\
&\quad + \Delta_s(h_2(V_s) - \tilde{h}_2(V_s)) - \int_0^{V_s} \rho_3(t, \mathbf{O}; \tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s) (h_2(V_s) - \tilde{h}_2(V_s)) d\tilde{\Lambda}_s(t), \tag{4.36}
\end{aligned}$$

where  $(\boldsymbol{\theta}^*, \Lambda_s^*)$  is equal to  $\varepsilon^*(\boldsymbol{\theta}, \Lambda_s) + (1 - \varepsilon^*)(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)$  for some  $\varepsilon^* \in [0, 1]$ . Thus, we have

$$\begin{aligned}
&|l_{\boldsymbol{\theta}}(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)^T \tilde{\mathbf{h}}_1 - l_{\Lambda_s}(\tilde{\boldsymbol{\theta}}, \tilde{\Lambda}_s)[\tilde{h}_2]| \\
&\leq e^{q_1 + q_2} \|\mathbf{Y}\| \left\{ \|\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}}\| + \|\mathbf{h}_1 - \tilde{\mathbf{h}}_1\| + |\Lambda_s(V_s) - \tilde{\Lambda}_s(V_s)| \right. \\
&\quad + \int_0^{\tau} |\Lambda_s(t) - \tilde{\Lambda}_s(t)| [dt + d|h_2(t)| + d|\tilde{h}_2(t)|] \\
&\quad \left. + |h_2(V_s) - \tilde{h}_2(V_s)| + \int_0^{\tau} |h_2(V_s) - \tilde{h}_2(V_s)| [d\Lambda_s(t) - d\tilde{\Lambda}_s(t)] \right\}, \tag{4.37}
\end{aligned}$$

where  $d|h_2(t)| = dh_2^+(t) + dh_2^-(t)$  and  $d|\tilde{h}_2(t)| = d\tilde{h}_2^+(t) + d\tilde{h}_2^-(t)$ . As done in Section 4.4.3.1, by applying Theorem 2.7.5 (p159) in van der Vaart and Wellner (1996), we note that for a set  $\mathcal{H}_2 = \{h_2 : \|h_2\|_V \leq B_1\}$ ,  $\log N_{[\cdot]}(\varepsilon, \mathcal{H}_2, L_2(P)) \leq K/\varepsilon$  for a constant  $B_1$  and any probability measure  $P$  where  $K$  is a constant. Thus, we can find  $\exp\{K/\varepsilon\}$  brackets,  $\{[L_j, U_j]\}$ , to cover the class  $\mathcal{H}_2$  such that  $\|U_j - L_j\|_{L_2(P)} \leq \varepsilon$  for each pair of  $[L_j, U_j]$ . On the other hand, we can further find a partition of  $\mathcal{H}_1 = \{\mathbf{h}_1 : \|\mathbf{h}_1\| \leq 1\}$ , say  $I_1 \cup I_2 \cup \dots$ , such that the number of partitions is of the order  $(1/\varepsilon)$ , and, for any  $\mathbf{h}_1$  and  $h_2$  in the same partition, their Euclidean distance is less than  $\varepsilon$ . Therefore, the

partition  $\{I_1, I_2, \dots\} \times \{[L_j, U_j]\}$  bracket covers  $\mathcal{H}_1 \times \mathcal{H}_2$ , and the total number of the partitions is of order  $(1/\varepsilon) \exp\{1/\varepsilon\}$ . Then, we obtain

$$\log N_{[\cdot]}(\varepsilon, \mathcal{G}, L_2(P)) \leq O\left(\frac{1}{\varepsilon} + \log \varepsilon\right).$$

Moreover,  $\mathcal{G}$  has an  $L_2(P)$ -integrable covering function, which is equal to  $O(e^{q_1+q_2\|\mathbf{Y}\|})$ .

Hence, from Theorem 2.5.6 (p130) in van der Vaart and Wellner (1996),  $\mathcal{G}$  is P-Donsker.

Additionally, from (4.37), we can calculate

$$\begin{aligned} & \left| l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right| \\ & \leq e^{q_1+q_2\|\mathbf{Y}\|} \left\{ \|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + |\Lambda_s(V_s) - \Lambda_{s0}(V_s)| + \int_0^\tau |\Lambda_s(t) - \Lambda_{s0}(t)| dt \right\} \\ & \quad + \left| \int_0^\tau \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d(\Lambda_s(t) - \Lambda_{s0}(t)) \right|. \end{aligned} \quad (4.38)$$

If  $\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| \rightarrow 0$  and  $\sup_{t \in [0, \tau]} |\Lambda_s(t) - \Lambda_{s0}(t)| \rightarrow 0$ , the above expression converges to zero uniformly. Therefore,

$$\sup_{(\mathbf{h}_1, h_2) \in \mathcal{H}} \mathbf{P} \left[ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right]^2 \longrightarrow 0.$$

#### 4.4.3.3 Derivative operator $\nabla U_{\xi_0}$

From (4.36) in the previous Section 4.4.3.2, we can obtain

$$\begin{aligned} & l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \\ & = [l_\theta(\boldsymbol{\theta}, \Lambda_s)^T - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T] \mathbf{h}_1 + [l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s) - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})][h_2] \\ & = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \nabla_{\boldsymbol{\theta}} \rho_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \mathbf{h}_1 - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d\Lambda_s^*(t) \\ & \quad + \int_0^{V_s} \rho_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \Lambda_{s0})(t) \\ & \quad + \int_0^{V_s} \int_0^t \rho_5(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \Lambda_{s0})(u) \mathbf{h}_1 d\Lambda_s^*(t) \end{aligned}$$

$$\begin{aligned}
& - \int_0^{V_s} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T \mathbf{h}_1 d(\Lambda_s - \Lambda_{s0})(t) \\
& - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
& + \int_0^{V_s} \int_0^t \rho_6(u, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d(\Lambda_s - \Lambda_{s0})(u) h_2(t) d\Lambda_s^*(t) \\
& - \int_0^{V_s} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T h_2(t) d(\Lambda_s - \Lambda_{s0})(t) \\
& = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \left\{ \nabla_{\boldsymbol{\theta}} \rho_1(\mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*)^T d\Lambda_s^*(t) \right\} \mathbf{h}_1 \\
& + \mathbf{h}_1^T \left\{ \int_0^\tau I(t \leq V_s) [\rho_4(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \rho_2(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \right. \\
& \quad \left. + \rho_5(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \int_t^{V_s} d\Lambda_s^*(u)] d(\Lambda_s - \Lambda_{s0})(t) \right\} \\
& - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^\tau I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) d\Lambda_s^*(t) \\
& - \int_0^\tau \left\{ -I(t \leq V_s) \rho_6(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) \int_t^{V_s} h_2(u) d\Lambda_s^*(u) \right. \\
& \quad \left. + I(t \leq V_s) \rho_3(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) h_2(t) \right\} d(\Lambda_s - \Lambda_{s0})(t). \tag{4.39}
\end{aligned}$$

Then, we have

$$\begin{aligned}
& \nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] \\
& = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \mathbf{P} \left\{ \nabla_{\boldsymbol{\theta}} \rho_1(\mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\Lambda_{s0}(t) \right\} \mathbf{h}_1 \\
& + \mathbf{h}_1^T \left\{ \int_0^\tau \mathbf{P} \left[ I(t \leq V_s) \left( \rho_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \rho_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right. \right. \right. \\
& \quad \left. \left. + \rho_5(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u) \right) \right] d(\Lambda_s - \Lambda_{s0})(t) \right\} \\
& - (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \int_0^\tau \mathbf{P} \left\{ I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} h_2(t) d\Lambda_{s0}(t) \\
& - \int_0^\tau \mathbf{P} \left\{ -I(t \leq V_s) \rho_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} h_2(u) d\Lambda_{s0}(u) \right. \\
& \quad \left. + I(t \leq V_s) \rho_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) h_2(t) \right\} d(\Lambda_s - \Lambda_{s0})(t).
\end{aligned}$$

By the similar algebra done in (4.38), we can verify that, for  $j = 1, \dots, 6$ ,

$$\sup_{t \in [0, \tau]} \|\rho_j(t, \mathbf{O}; \boldsymbol{\theta}^*, \Lambda_s^*) - \rho_j(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\| \leq e^{q_3 + q_4 \|\mathbf{Y}\|} \left\{ \|\boldsymbol{\theta}^* - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s^* - \Lambda_{s0}| \right\},$$

which implies that the linear operator  $\nabla U_{\xi_0}$  is bounded.

Then, we obtain

$$\begin{aligned} & \mathbf{P} \left[ l_\theta(\boldsymbol{\theta}, \Lambda_s)^T \mathbf{h}_1 + l_{\Lambda_s}(\boldsymbol{\theta}, \Lambda_s)[h_2] - l_\theta(\boldsymbol{\theta}_0, \Lambda_{s0})^T \mathbf{h}_1 - l_{\Lambda_s}(\boldsymbol{\theta}_0, \Lambda_{s0})[h_2] \right] \\ &= \nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] + o\left(\|\boldsymbol{\theta} - \boldsymbol{\theta}_0\| + \sup_{t \in [0, \tau]} |\Lambda_s - \Lambda_{s0}|\right)(\|\mathbf{h}_1\| + \|h_2\|_V). \end{aligned}$$

Therefore,  $U_\xi$  is *Fréchet* differentiable at  $\xi_0$ .

Additionally, from (4.39) and the above expression, we have

$$\nabla U_{\xi_0}(\boldsymbol{\theta} - \boldsymbol{\theta}_0, \Lambda_s - \Lambda_{s0})[\mathbf{h}_1, h_2] = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^T \Omega_1[\mathbf{h}_1, h_2] + \int_0^\tau \Omega_2[\mathbf{h}_1, h_2] d(\Lambda_s - \Lambda_{s0})(t),$$

where

$$\begin{aligned} \Omega_1[\mathbf{h}_1, h_2] &= \mathbf{P} \left\{ \nabla_{\boldsymbol{\theta}} \rho_1(\mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \int_0^{V_s} \nabla_{\boldsymbol{\theta}} \rho_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) d\Lambda_{s0}(t) \right\} \mathbf{h}_1 \\ &\quad - \int_0^\tau \mathbf{P} \left\{ I(t \leq V_s) \nabla_{\boldsymbol{\theta}} \rho_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} h_2(t) d\Lambda_{s0}(t) \end{aligned}$$

and

$$\begin{aligned} & \Omega_2[\mathbf{h}_1, h_2] \\ &= \mathbf{h}_1^T \mathbf{P} \left\{ I(t \leq V_s) [\rho_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \rho_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) + \rho_5(u, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u)] \right\} \\ &\quad + \mathbf{P} \left\{ I(t \leq V_s) \rho_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} h_2(u) d\Lambda_{s0}(u) \right\} \\ &\quad - \mathbf{P} \left\{ I(t \leq V_s) \rho_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right\} h_2(t). \end{aligned}$$



Thus,  $\Omega = (\Omega_1, \Omega_2)$  is the bounded linear operator from  $R^d \times BV[0, \tau]$  to itself. Furthermore, we note that  $\Omega = \mathbf{H} + (\mathbf{M}_1, \mathbf{M}_2)$ , where

$$\begin{aligned}\mathbf{H}(\mathbf{h}_1, h_2) &= (\mathbf{h}_1, -\mathbf{P}\{I(t \leq V_s)\rho_3(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0})\}h_2(t)), \\ \mathbf{M}_1(\mathbf{h}_1, h_2) &= \Omega_1[\mathbf{h}_1, h_2] - \mathbf{h}_1, \\ \mathbf{M}_2(\mathbf{h}_1, h_2) &= \mathbf{h}_1^T \mathbf{P} \left\{ I(t \leq V_s) \left[ \nabla_{\boldsymbol{\theta}} \rho_4(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) - \rho_2(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \right. \right. \\ &\quad \left. \left. + \rho_5(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} d\Lambda_{s0}(u) \right] \right\} \\ &\quad + \mathbf{P} \left\{ I(t \leq V_s) \rho_6(t, \mathbf{O}; \boldsymbol{\theta}_0, \Lambda_{s0}) \int_t^{V_s} h_2(u) d\Lambda_{s0}(u) \right\},\end{aligned}$$

and also note that  $\mathbf{H}$  is obviously invertible. Since  $\mathbf{M}_1$  maps into a finite-dimensional space, it is compact. The image of  $\mathbf{M}_2$  is a continuously differentiable function in  $[0, \tau]$ . By the Arzela-Ascoli theorem (p41) in van der Vaart and Wellner (1996),  $\mathbf{M}_2$  is a compact operator from  $R^d \times BV[0, \tau]$  to  $BV[0, \tau]$ . Thus, we conclude that  $\Omega$  is the summation of an invertible operator  $\mathbf{H}$  and a compact operator  $\mathbf{M} = (\mathbf{M}_1, \mathbf{M}_2)$ .

## 4.5 Simulation Studies

In this section, we present the results from our simulation studies. First, to assess finite sample properties of the proposed maximum likelihood estimators, two sets of simulations with different generalized linear mixed models for the longitudinal outcomes are performed. Continuous and binary data are considered for longitudinal process in the simulations in Sections 4.5.1 and 4.5.2, respectively. Then, we conduct simulation studies for robustness of the assumed mixture distribution in Section 4.5.3. Selection procedures for the number of mixtures by AIC and BIC criteria are assessed through simulation studies in Section 4.5.4.

### 4.5.1 Continuous longitudinal outcomes and survival time

In this section, we assume  $Y_{ij}$  follows a Gaussian distribution given a subject-specific random intercept. Specifically we have

$$Y_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + b_i + \epsilon_{ij} = \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3ij} + b_i + \epsilon_{ij},$$

for  $j = 1, \dots, n_i$ , where  $\epsilon_{ij} \sim \mathcal{N}(0, \sigma_y^2)$ , and

$$h(t|b_i) = \lambda(t) \exp\{\psi b_i + \mathbf{Z}_i(t)\boldsymbol{\gamma}\} = \lambda(t) \exp\{\psi b_i + \gamma_1 Z_{1i} + \gamma_2 Z_{2i}\},$$

where  $b_i \sim \sum_{k=1}^K w_k \mathcal{N}(\mu_k, \sigma_b^2)$ ,  $K$  is the number of mixture components, and  $K = 2$  and  $K = 3$  are simulated.  $X_{1i} \equiv Z_{1i}$  are simulated from a Bernoulli distribution with success probability being 0.5, and  $X_{2i} \equiv Z_{2i}$  are simulated from the uniform distribution between 0 and 1. The longitudinal data are generated for every 0.1 unit of time, and thus  $X_{3ij}$ , the time at measurement, has the value of every 0.1 unit ranging over 0 through 2.4. We consider  $\psi = -0.1$  indicating negative dependency between longitudinal process and survival time model. The parameters in the longitudinal and hazard models are chosen as  $\beta_1 = 1$ ,  $\beta_2 = -0.5$ ,  $\beta_3 = -0.2$ ,  $\sigma_y^2 = 0.5$ ,  $\psi = -0.1$ ,  $\gamma_1 = -0.1$ ,  $\gamma_2 = 0.1$ , and  $\lambda(t) = 1$ . The parameters in the mixture distribution for random effects are  $\mu_1 = -1.5$ ,  $\mu_2 = 1.5$ , and  $w_1 = 0.4$  for  $K = 2$  and  $\mu_1 = -3$ ,  $\mu_2 = 0$ ,  $\mu_3 = 3$ ,  $w_1 = 0.4$ , and  $w_2 = 0.3$  for  $K = 3$ . The weight of the last mixture component ( $w_2$  and  $w_3$  for  $K = 2$  and  $K = 3$  respectively) is determined from the restriction  $\sum_{k=1}^K w_k = 1$ . The variance of random effects  $\sigma_b^2$  is chosen as 0.3. Censoring time is generated from the uniform distribution between 0.4 and 2.4, and the censoring proportion is around 25~35%. We consider different sample sizes ( $n=400, 800$ ) with 1000 replications. The average number of longitudinal observations ( $n_i$ ) is 7–8 with the range of 1 to 24. For the estimated baseline cumulative hazard function, we consider three fixed time points of 0.9, 1.4, and 1.9.

The results of the maximum likelihood estimates for  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \sigma_y^2, \boldsymbol{\mu}^T, \boldsymbol{w}^T, \sigma_b^2, \psi, \boldsymbol{\gamma}^T)^T$  and baseline cumulative hazards at the three time points and their respective standard error estimates are reported in Table 4.1. In Table 4.1, “True” gives the true values of parameters; the averages of the maximum likelihood estimates from the EM algorithm are in “Est.”; the sample standard deviations from 1000 simulations are reported in “SSD”; “ESE” is the average of 1000 standard error estimates based on the observed information matrix; “CP” is the coverage proportion of 95% confidence intervals based on the estimated standard error “ESE”. Satterthwaite method is used for the coverage probabilities of  $\sigma_y^2$  and  $\sigma_b^2$ .

From Table 4.1, we can see that even for the smaller sample size (n=400), the bias of the estimates from EM algorithm is negligible for most cases. The estimated standard errors calculated from the observed information matrix are close to the sample standard deviations from the 1000 estimates, and the 95% confidence interval coverage rates are close to 0.95 except for weights of the mixture components. The coverage rates of weights are improved for larger sample size in both 2 and 3 mixtures. The estimates for the parameters in the longitudinal and hazards models ( $\boldsymbol{\beta}$ ,  $\sigma_y^2$ ,  $\psi$ ,  $\boldsymbol{\gamma}$  and  $\Lambda(t)$ ) perform well for different mixtures.

## 4.5.2 Binary longitudinal outcomes and survival time

In this section, we assume that  $Y_{ij}$  is a binary outcome following

$$P(Y_{ij} = y_{ij}|b_i) = \exp \left\{ y_{ij}\eta_{ij} - \log(1 + \exp\{\eta_{ij}\}) \right\}, \quad y_{ij} = 0, 1,$$

with  $\eta_{ij} = \boldsymbol{X}_{ij}\boldsymbol{\beta} + b_i = \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3ij} + b_i$  for  $j = 1, \dots, n_i$ , and we consider the same hazards model and simulation setting as those used in Section 4.5.1 except the followings. The parameters in the mixture distribution for random effects are  $\mu_1 = -3$ ,

Table 4.1: Summary of simulation results of maximum likelihood estimation using mixtures of Gaussian distributions for random effects in the joint modeling of continuous longitudinal outcomes and survival time.

Mixture	Par.	True	n=400				n=800			
			Est.	SSD	ESE	CP	Est.	SSD	ESE	CP
2	$\beta_1$	1.0	.983	.066	.068	.958	.985	.047	.048	.947
	$\beta_2$	- .5	- .529	.107	.119	.969	- .540	.079	.084	.947
	$\beta_3$	- .2	- .203	.033	.033	.955	- .203	.024	.024	.952
	$\sigma_y^2$	.5	.500	.014	.014	.954	.500	.010	.010	.948
	$\mu_1$	-1.5	-1.478	.081	.088	.962	-1.469	.060	.062	.938
	$\mu_2$	1.5	1.524	.075	.082	.966	1.530	.055	.058	.940
	$w_1$	.4	.400	.025	.033	.991	.401	.018	.023	.981
	$\sigma_b^2$	.3	.296	.029	.029	.955	.298	.020	.020	.958
	$\psi$	- .1	- .102	.040	.039	.950	- .100	.028	.028	.946
	$\gamma_1$	- .1	- .101	.123	.121	.945	- .105	.085	.085	.952
	$\gamma_2$	.1	.102	.209	.210	.954	.096	.144	.147	.950
	$\Lambda(.9)$	.9	.911	.130	.128	.950	.909	.087	.090	.955
	$\Lambda(1.4)$	1.4	1.421	.206	.202	.942	1.415	.139	.141	.952
	$\Lambda(1.9)$	1.9	1.939	.304	.295	.953	1.924	.205	.205	.950
3	$\beta_1$	1.0	.983	.070	.071	.947	.984	.049	.050	.956
	$\beta_2$	- .5	- .543	.116	.123	.952	- .543	.085	.087	.922
	$\beta_3$	- .2	- .203	.034	.034	.949	- .204	.024	.024	.960
	$\sigma_y^2$	.5	.500	.014	.014	.957	.500	.010	.010	.950
	$\mu_1$	-3.0	-2.970	.084	.090	.954	-2.968	.064	.063	.909
	$\mu_2$	.0	.028	.093	.097	.954	.032	.069	.068	.933
	$\mu_3$	3.0	3.030	.089	.094	.954	3.034	.063	.066	.925
	$w_1$	.4	.400	.025	.033	.992	.400	.018	.023	.983
	$w_2$	.3	.299	.024	.029	.980	.300	.017	.020	.977
	$\sigma_b^2$	.3	.295	.029	.029	.956	.298	.021	.021	.946
	$\psi$	- .1	- .101	.024	.024	.956	- .101	.017	.017	.941
	$\gamma_1$	- .1	- .091	.112	.119	.963	- .096	.085	.084	.950
	$\gamma_2$	.1	.088	.215	.207	.946	.114	.146	.146	.944
	$\Lambda(.9)$	.9	.913	.125	.127	.948	.897	.088	.088	.951
	$\Lambda(1.4)$	1.4	1.417	.202	.200	.949	1.402	.141	.140	.949
	$\Lambda(1.9)$	1.9	1.928	.297	.292	.946	1.908	.206	.204	.948

$\mu_2 = 3$ , and  $w_1 = 0.4$  for  $K = 2$  and  $\mu_1 = -6$ ,  $\mu_2 = 0$ ,  $\mu_3 = 6$ ,  $w_1 = 0.4$ , and  $w_2 = 0.3$  for  $K = 3$ . The binary longitudinal data are generated for every 0.1 and 0.05 units of time for the mixture of 2 and 3 normal distributions, respectively, and  $X_{3ij}$ , the time at measurement, has the values of every 0.1 and 0.05 units corresponding to the mixture distributions ranging over 0 through 2.4. Thus, the average numbers of longitudinal observations ( $n_i$ ) are 7–8 with the range of 1 to 24 and 15–16 with the range of 1 to 48 for the mixture of 2 and 3 distributions, respectively.

The results of the maximum likelihood estimates for  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \boldsymbol{\mu}^T, \boldsymbol{w}^T, \sigma_b^2, \psi, \boldsymbol{\gamma}^T)^T$  and baseline cumulative hazards at the given three time points and their respective standard error estimates are reported in Table 4.2. Similar to the results for the continuous longitudinal outcomes, Table 4.2 shows that overall the estimates perform well even for the smaller sample size  $n = 400$  with small biases. The parameters of interest in longitudinal and hazards models have the estimated standard errors which are close to the sample standard deviations. Meanwhile, the estimated standard errors of the parameters of mixture components which are means of random effects and weights appear to be overestimated being larger than their sample standard deviations, which leads to the wide confidence interval.

### 4.5.3 Sensitivity for model-misspecification

In this section, we conduct simulation studies to examine the sensitivity of the assumed mixture distribution. We consider continuous longitudinal outcomes and survival time with the same setting used in Section 4.5.1 except for the true distribution of random effects. Random effects are generated from a mixture of a t-distribution with 10 degrees of freedom and non-centrality of -1 and a Gamma distribution with shape and scale parameters of 7 and 1/8 respectively. We assume equal probability for the two distributions. We fit 5 sets of simultaneous models assuming different mixtures for random

Table 4.2: Summary of simulation results of maximum likelihood estimation using mixtures of Gaussian distributions for random effects in the joint modeling of binary longitudinal outcomes and survival time.

Mixture	Par.	True	n=400				n=800			
			Est.	SSD	ESE	CP	Est.	SSD	ESE	CP
2	$\beta_1$	1.0	1.029	.193	.201	.960	1.015	.143	.141	.942
	$\beta_2$	- .5	- .508	.292	.323	.966	- .495	.205	.227	.965
	$\beta_3$	- .2	- .200	.166	.180	.966	- .203	.116	.127	.968
	$\mu_1$	-3.0	-3.046	.241	.275	.968	-3.034	.164	.193	.970
	$\mu_2$	3.0	3.016	.211	.253	.976	3.011	.142	.177	.984
	$w_1$	.4	.401	.025	.033	.993	.400	.017	.023	.991
	$\sigma_b^2$	.3	.329	.133	.195	.940	.332	.092	.136	.956
	$\psi$	- .1	- .099	.021	.021	.949	- .099	.015	.015	.955
	$\gamma_1$	- .1	- .103	.121	.122	.959	- .098	.087	.086	.947
	$\gamma_2$	.1	.091	.210	.211	.944	.104	.142	.149	.958
	$\Lambda(.9)$	.9	.910	.131	.130	.955	.900	.088	.091	.956
	$\Lambda(1.4)$	1.4	1.421	.209	.206	.934	1.402	.142	.143	.956
	$\Lambda(1.9)$	1.9	1.932	.310	.299	.941	1.899	.205	.207	.948
3	$\beta_1$	1.0	.988	.167	.171	.953	.993	.123	.121	.947
	$\beta_2$	- .5	- .519	.268	.287	.960	- .516	.189	.203	.967
	$\beta_3$	- .2	- .208	.126	.128	.957	- .206	.091	.091	.951
	$\mu_1$	-6.0	-5.844	.353	.483	.967	-5.872	.260	.342	.963
	$\mu_2$	.0	.023	.172	.194	.970	.018	.127	.138	.966
	$\mu_3$	6.0	6.024	.397	.504	.984	6.006	.303	.349	.971
	$w_1$	.4	.402	.025	.035	.995	.402	.018	.024	.989
	$w_2$	.3	.298	.025	.034	.986	.298	.017	.024	.985
	$\sigma_b^2$	.3	.277	.095	.100	.977	.289	.070	.072	.966
	$\psi$	- .1	- .102	.014	.015	.955	- .101	.011	.010	.946
	$\gamma_1$	- .1	- .103	.121	.120	.955	- .107	.085	.084	.948
	$\gamma_2$	.1	.104	.201	.208	.961	.099	.147	.146	.949
	$\Lambda(.9)$	.9	.909	.128	.130	.950	.911	.094	.092	.930
	$\Lambda(1.4)$	1.4	1.421	.202	.207	.960	1.420	.147	.146	.946
	$\Lambda(1.9)$	1.9	1.926	.297	.302	.958	1.929	.220	.213	.946

effects which are 1 normal distribution without mixture and the mixtures of 2, 3, 4 and 5 normal distributions, and we compare the results for the parameters of interest in longitudinal and hazards models and the estimated density plots of random effects. Table 4.3 shows the results of longitudinal and hazards models from assuming the 5 different models for random effects. We can see that bias gets smaller and coverage rates become closer to the 95% nominal level as the number of mixtures increases. From the table, we also find that more mixture produces estimates more close to the true values in the longitudinal model while estimates in hazards model are less sensitive to the number of distributions in mixture. In other words, when the true distribution of random effects is not a Gaussian distribution, the use of mixture is effective in longitudinal model but the inference on hazards model is reasonable regardless of mixture. Figure 4.1 shows the true and estimated density plots of random effects. From these density plots, all the mixture models of 2, 3, 4 and 5 normal distributions produces similar shapes to the true distribution while one normal distribution does not. The mixture of 5 normal distribution appears to be close to the true density. Figure 4.2 shows the relative bias plot of the parameters in longitudinal and hazard models which are denoted with thin and thick lines respectively. The relative biases are calculated from the median absolute biases divided by their absolute true values. This Figure 4.2 confirms what we observe in Table 4.3.

#### **4.5.4 Selection of the number of mixture distributions**

We adopt Akaike Information Criterion (AIC) and Bayesian information criterion (BIC) for selection of the number of normal distribution in mixture and assess these selection procedures through simulation studies in this section. AIC gives a penalty to a model with more parameters and BIC gives a penalty to a model with more parameters and larger sample size. Given a data set, competing models are ranked according to their

Table 4.3: Summary of simulation results of sensitivity for model-misspecification

Par.	1 Normal distribution					2 Normal distributions					3 Normal distributions					4 Normal distributions					5 Normal distributions						
	TRUE	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP	Est.	SSD	ESE	CP		
< Longitudinal model >																											
$\beta_1$	1.0	.947	.085	.095	.931	.965	.074	.074	.928	.977	.061	.060	.935	.978	.061	.060	.934	.985	.055	.054	.943	.985	.055	.054	.943		
$\beta_2$	.5	.591	.122	.164	.975	.582	.118	.128	.933	.557	.101	.105	.930	.555	.101	.104	.928	.536	.092	.095	.928	.536	.092	.095	.928		
$\beta_3$	.2	.203	.024	.025	.964	.202	.024	.025	.961	.202	.024	.024	.961	.202	.024	.024	.959	.202	.024	.024	.955	.202	.024	.024	.955		
$\sigma_y^2$	.5	.501	.010	.010	.948	.501	.010	.010	.948	.501	.010	.010	.950	.501	.010	.010	.952	.500	.010	.010	.946	.500	.010	.010	.946		
< Hazards model >																											
$\psi$	.1	.102	.034	.033	.942	.101	.034	.033	.943	.101	.034	.033	.941	.101	.034	.033	.940	.102	.034	.033	.945	.102	.034	.033	.945		
$\gamma_1$	.1	.093	.083	.085	.949	.095	.083	.084	.953	.096	.083	.084	.954	.096	.083	.084	.951	.098	.085	.084	.948	.098	.085	.084	.948		
$\gamma_2$	.1	.107	.144	.147	.949	.105	.144	.147	.949	.103	.143	.146	.951	.102	.143	.146	.951	.108	.144	.146	.952	.108	.144	.146	.952		
$\Lambda(.9)$	.9	.906	.089	.089	.943	.907	.089	.089	.944	.906	.089	.089	.945	.907	.089	.089	.943	.904	.090	.089	.944	.904	.090	.089	.944		
$\Lambda(1.4)$	1.4	1.408	.139	.140	.945	1.409	.139	.140	.945	1.408	.139	.140	.947	1.409	.139	.140	.946	1.403	.137	.139	.947	1.403	.137	.139	.947		
$\Lambda(1.9)$	1.9	1.911	.202	.203	.956	1.911	.202	.203	.956	1.910	.200	.203	.958	1.911	.202	.203	.956	1.905	.200	.202	.954	1.905	.200	.202	.954		



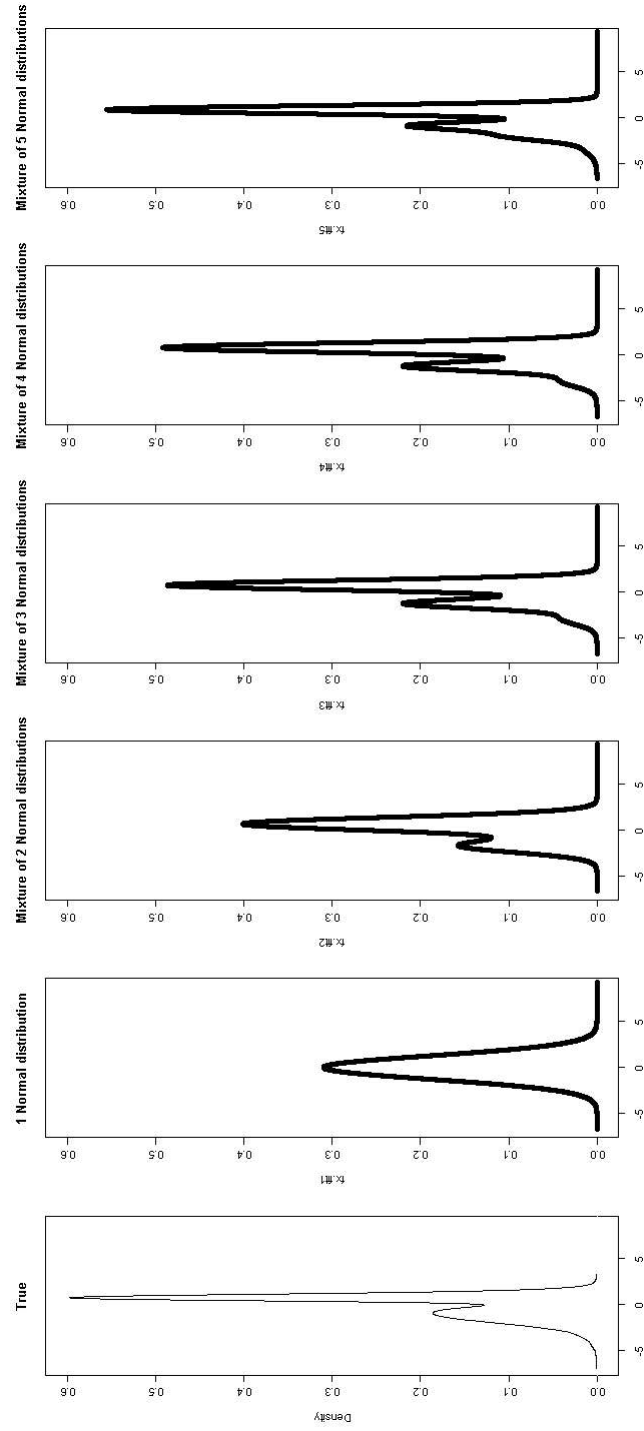


Figure 4.1: Density plots of random effects from simulation results of sensitivity for model-misspecification

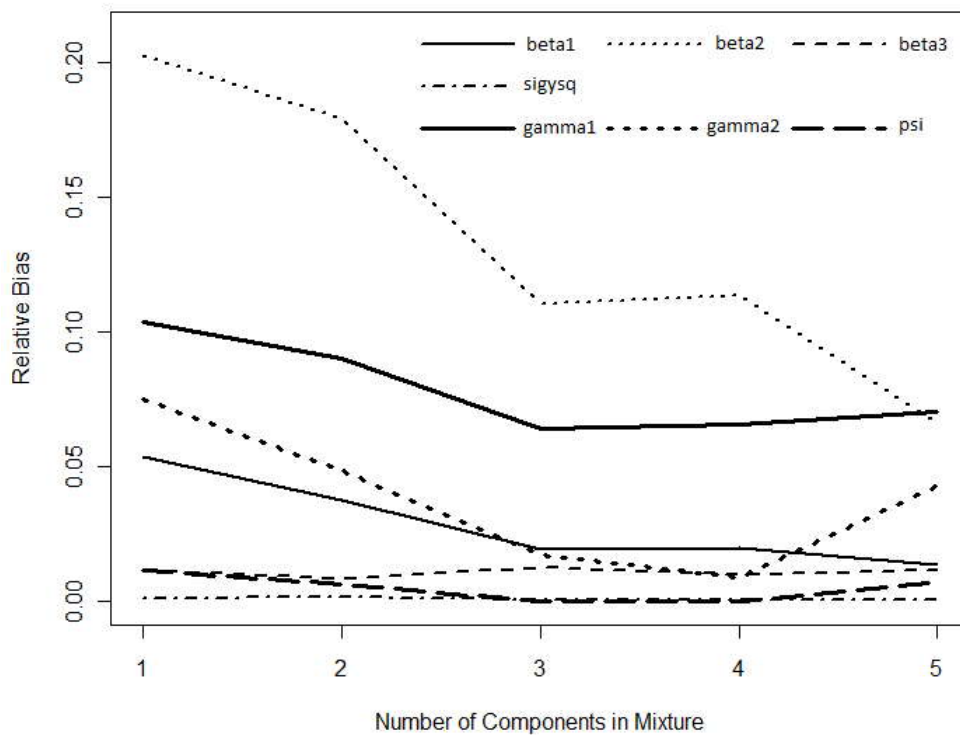


Figure 4.2: Relative bias plot of parameters in longitudinal and hazard models (thin and thick lines respectively) from simulation results of sensitivity for model-misspecification

Table 4.4: Summary of simulation results: Frequencies on the selected number of Normal distributions in mixture (n=200)

Criteria	Selected # of Normal distributions in mixture				
	1	2	3	4	5
AIC	0	0	969	21	10
BIC	0	0	990	5	5

AIC (or BIC), with the one having the lowest AIC (or BIC) being the best. Continuous longitudinal outcomes and survival time are considered with the same setting used in Section 4.5.1. Random effects are generated from a mixture of 3 normal distributions. We fit 5 sets of simultaneous models with different mixtures for random effects which are 1 normal distribution without mixture and the mixtures of 2, 3, 4 and 5 normal distributions. AIC and BIC values are calculated for all 5 fitted mixture models in each data set and we report frequencies of mixture models selected as best by AIC and BIC among 1000 data sets. We consider sample sizes of 200 and 800.

In Table 4.4, we summarize the results for the sample size of 200. We see that both AIC and BIC mostly select the true distribution of a mixture of 3 normal distributions as best. For the large sample size of 800, the mixture of 3 normal distributions is selected by both AIC and BIC for all 1000 simulated data sets. This demonstrates that the number of mixture distributions is properly selected by AIC and BIC even for small sample sizes.

## 4.6 Analysis of the CHANCE Study

The Carolina Head and Neck Cancer Study (CHANCE) is a population based epidemiologic study conducted at 60 hospitals in 46 counties in North Carolina from 2002 through 2006 (Divaris *et al.* 2010). Patients were diagnosed with head and neck cancer (oral, pharynx, and larynx cancer) from 2002–2006. Their survival status was collected

up to 2007 and QoL was evaluated over time for three years after diagnosis. QoL information was collected through questionnaires. Based on summary scores of the five domains of self-perceived quality of life including Physical Well-Being (PWB), Social/Family Well-Being (SWB), Emotional Well-Being (EWB), Functional Well-Being (FWB) and Head and Neck Cancer Specific symptoms (HNCS), patient's QoL information was classified into satisfaction or dissatisfaction with life. Survival time is defined as the time to death from diagnosis. Demographic and life style characteristics, medical histories and clinical factors are also collected. Ending in December 2009 and excluding the patients with missing data, information on QoL has been obtained from 554 head and neck cancer patients. Based on the death information through 2007 available from the National Death Index (NDI), 85 of 554 patients died and the censoring rate is 85%. The number of observations per patient ranges 1 to 3 with average of 1.93. It is of interest to elucidate the variables which are associated with both QoL satisfaction and survival time for patients with head and neck cancer. In particular, we are interested in the comparison between African-Americans and Whites since it is known that African-Americans have a higher incidence of head and neck cancer and worse survival than Whites. The longitudinal QoL satisfaction outcomes and survival time are correlated within a patient, and this dependency should be taken into account in the analysis.

We apply our proposed method to Head and Neck Cancer Specific symptoms (HNCS) among QoL domains with survival time. Longitudinal HNCS QoL outcomes are binary measurements with 1 ("satisfied") and 0 ("dissatisfied"). We are interested in investigating which factors are related to QoL satisfaction and the risk of death. In the full models for both longitudinal QoL and survival time, we consider race (African-Americans, Whites), the number of 12 oz. beers consumed per week (None, <1, 1–4, 5–14, 15–29,  $\geq 30$ ), household income (0–10K, 20–30K, 40–50K,  $\leq 60$ K), surgery (Yes/No), radiation therapy (Yes/No), chemotherapy (Yes/No), primary tumor site

(Oral & Pharyngeal, Laryngeal) and tumor stage (I, II, III, IV) as categorical, and age at diagnosis (range: 24–80), the number of persons supported by household income (range: 1–5), body mass index (BMI) (range: 15.66–56.28) and the total number of medical conditions reported (range: 0–6) as continuous. Additionally, 2 interactions with race, i.e. race  $\times$  the total number of medical conditions reported and race  $\times$  tumor site, are included in both models since we are particularly interested in the difference of QoL and survival between African American and White. Time at survey measurement is also included as a covariate for longitudinal outcomes. A random intercept for the dependence between the QoL satisfaction and the risk of death is included in both models, and assumed to follow a mixture of normal distributions.

For the full model, we first considered 5 different distributions for random effects which are 1 normal distribution without mixture and the mixtures of 2, 3, 4 and 5 normal distributions, and both AIC and BIC selected a mixture of 3 normal distributions with their lowest values as best. Then, we conducted backward variable selection based on the Likelihood Ratio Test (LRT) from the full model assuming the mixture of 3 normal distributions for random effects. Table 4.5 gives the results from the final models after removing non-significant covariates by LRT. From the “Simultaneous” columns, we see the number of 12 oz. beers consumed per week, household income and tumor stage are significantly associated with both patients’ HNCS QoL satisfaction and hazard of death. Using 30 or more of 12 oz. beers consumed per week as the reference group, all categories of the smaller amount are associated with higher odds of being satisfied while the categories of ‘none’ and ‘5 to 14’ of 12 oz. beers consumed per week are associated with lower risk of death. Higher household income is generally associated with higher odds of being satisfied and lower risk of death. Both patients’ HNCS QoL satisfaction and risk of death are significantly different for patients in different tumor stages. On the other hand, race (African-American), radiation therapy, the number of

Table 4.5: Results from final models of simultaneous and separate analyses for the Quality of Life and survival time for the CHANCE study

Parameter		Simultaneous			Separate		
		Est.	ESE	P-value	Est.	ESE	P-value
<i>HNCS QoL longitudinal model</i>							
Intercept	$\beta_0$				1.190	.390	.002
Race (ref= White): African American	$\beta_1$	.900	.399	.024	.511	.256	.047
# of 12 oz. beers consumed per week (ref=30 or more)							
– None	$\beta_2$	.858	.428	.045	.622	.300	.038
– less than 1	$\beta_3$	1.119	.600	.062	.735	.396	.064
– 1 to 4	$\beta_4$	1.588	.563	.005	1.268	.326	<.001
– 5 to 14	$\beta_5$	1.450	.428	.001	1.018	.279	<.001
– 15 to 29	$\beta_6$	1.007	.531	.058	.547	.327	.095
Household income (ref= level1: 0–10K)							
– level2: 20–30K	$\beta_7$	– .337	.358	.346	– .328	.258	.204
– level3: 40–50K	$\beta_8$	.633	.440	.151	.250	.282	.376
– level4: $\geq$ 60K	$\beta_9$	1.960	.509	<.001	1.045	.286	<.001
Radiation therapy (ref= No) : Yes	$\beta_{10}$	–1.668	.608	.006	–1.048	.280	<.001
Tumor stage (ref= I)							
– II	$\beta_{11}$	– .683	.554	.218	– .352	.330	.286
– III	$\beta_{12}$	–2.012	.534	<.001	–1.198	.314	<.001
– IV	$\beta_{13}$	–1.826	.507	<.001	–1.057	.277	<.001
# of persons supported by household income	$\beta_{14}$	– .388	.140	.006			
BMI	$\beta_{15}$	.061	.026	.021			
Time at survey measurement (years)	$\beta_{16}$	.354	.093	<.001	.254	.067	<.001
<i>Hazards model</i>							
Random effect coefficient	$\psi$	– .206	.078	.008			
# of 12 oz. beers consumed per week (ref=30 or more)							
– None	$\gamma_1$	– .705	.347	.042			
– less than 1	$\gamma_2$	– .156	.393	.692			
– 1 to 4	$\gamma_3$	– .712	.385	.064			
– 5 to 14	$\gamma_4$	– .991	.348	.004			
– 15 to 29	$\gamma_5$	– .579	.370	.117			
Household income (ref= level1: 0–10K)							
– level2: 20–30K	$\gamma_6$	– .206	.274	.453	– .219	.263	.406
– level3: 40–50K	$\gamma_7$	– .884	.341	.010	– .928	.331	.005
– level4: $\geq$ 60K	$\gamma_8$	–1.401	.374	<.001	–1.393	.358	<.001
Tumor stage (ref= I)							
– II	$\gamma_9$	– .255	.443	.564	– .295	.435	.498
– III	$\gamma_{10}$	.168	.403	.677	.136	.389	.727
– IV	$\gamma_{11}$	.950	.306	.002	.914	.295	.002
Total # of medical conditions reported	$\gamma_{12}$	.207	.095	.030	.205	.091	.025

P-value for testing  $\sigma_b^2$  being zero is based on a mixture of 0 and  $\chi^2$  distribution with 1 degree of freedom with equal mixing probabilities.

persons supported by household income and BMI are selected only in the HNCS QoL longitudinal model while the number of medical conditions reported is significant only in the hazard model. The results indicate that African-Americans, patients not treated with radiation therapy, patients in the family with the smaller number of persons supported by household income, or patients with higher BMI are associated with higher odds of being satisfied, but the risk of death is not affected by these factors. On the other hand, higher number of reported medical conditions is associated with higher risk of death, but it is not associated with HNCS QoL satisfaction. Furthermore, time at survey measurement is statistically significant in the HNCS QoL longitudinal model implying that patients have higher odds to be satisfied over time. The parameter  $\psi$  for the dependence between longitudinal HNCS QoL and survival time is negative and is statistically significant with p-value as 0.008. This means the longitudinal HNCS QoL and survival time are correlated and some latent factors which increase HNCS QoL satisfaction also decrease the risk of death. Although not provided in Table 4.5, we have additional parameters of the mixture distribution for random effects in the simultaneous modeling. The obtained estimates of three means of random effects are -3.146, 0.376, and 1.730 with estimated standard errors of 1.284, 2.897 and 0.986 and p-values of 0.014, 0.897 and 0.079, respectively. The first and second mixture components have the weight estimates of 0.147 and 0.105 with estimated standard errors of 0.062 and 0.051 and p-values of 0.018 and 0.037, respectively, and the common variance estimate of random effects is 0.637 with its estimated standard error of 1.286 and p-value of 0.483. In particular, the two weights of mixture components are significant at significant level 0.05, which strengthens the mixture of 3 normal distributions with the estimated 3 means of random effects.

For the purpose of comparison, we also conducted separate analyses for longitudinal HNCS QoL and survival time whose results are given in the last three columns of Table

4.5. Comparing the results from the simultaneous and separate analyses in Table 4.5, we can see our simultaneous analysis identifies two additional factors (the number of persons supported by household income and BMI) in the HNCS QoL longitudinal model and one additional factor (the number of 12 oz. beers consumed per week) in the hazard model.

Figure 4.3 shows the estimated baseline cumulative hazard rates over follow-up time with 95% confidence interval. The estimated baseline cumulative hazard rates look flat at the very early time within a year, but soon appear to be linearly increasing. Figure 4.4 shows the predicted conditional longitudinal trend of HNCS QoL satisfaction probabilities based on the simultaneous models (solid line) and the empirical longitudinal trend of HNCS QoL satisfaction probabilities (dotted line) based on the empirical longitudinal HNCS QoL satisfaction probabilities (dots). The predicted conditional probability of HNCS QoL satisfaction is calculated as the conditional expectation of the conditional probability of HNCS QoL satisfaction given the subject is alive at time  $t$ . That is,  $E_{b,\alpha} [P(Y(t) = 1|T > t) | \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}]$  using model notations in Section 4.2. The empirical probability of HNCS QoL satisfaction is calculated for every 0.05 unit of time at survey measurements. From Figure 4.4, the longitudinal trend of HNCS QoL satisfaction probabilities appears to be increasing over time and the empirical probabilities also gradually increase over time.

## 4.7 Concluding Remarks

We have relaxed normality assumption of random effects in the simultaneous modeling of longitudinal outcomes and survival time. Assuming the underlying distribution of random effects to be unknown, we used a mixture of Gaussian distributions as an approximation for the random effect distribution. We developed a maximum likelihood estimation method for the proposed simultaneous models and presented asymptotic



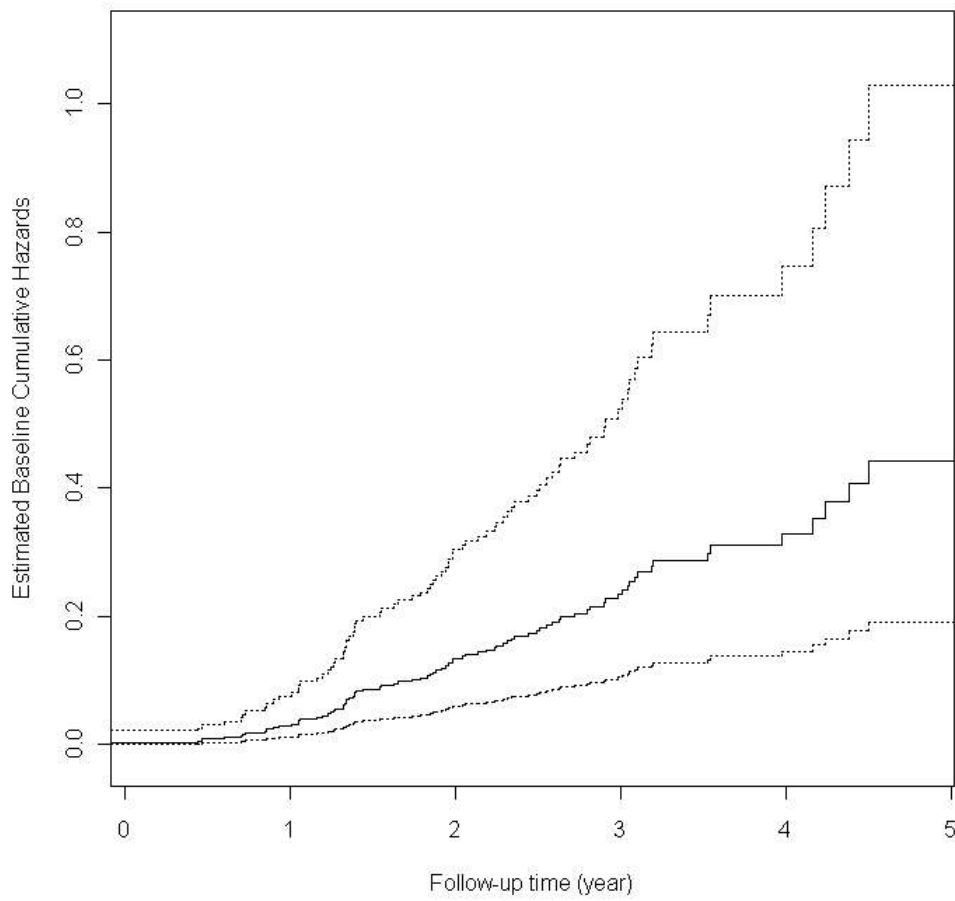


Figure 4.3: Estimated baseline cumulative hazards (solid line) with 95% confidence interval (dotted lines) by the simultaneous analysis of HNCS QoL longitudinal outcome and survival time

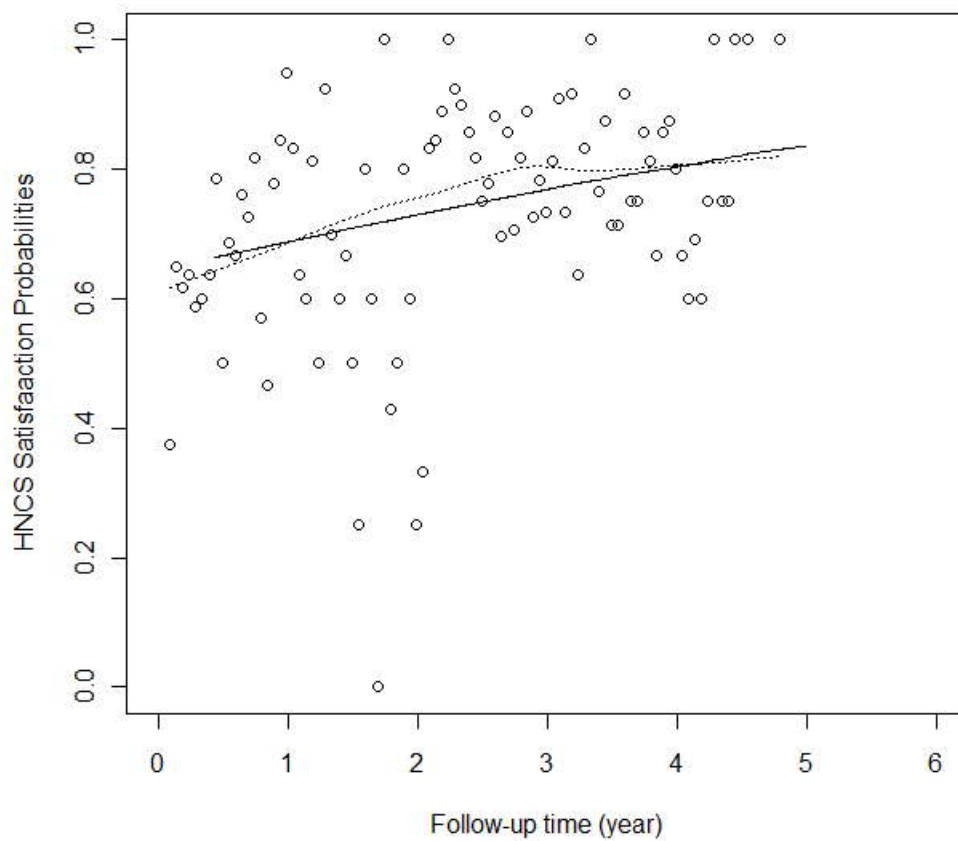


Figure 4.4: The predicted conditional longitudinal trend based on the simultaneous models (solid line) and the empirical longitudinal trend (dotted line) based on the empirical longitudinal HNCS QoL satisfaction probabilities (dots)

properties of the proposed estimators. The proposed estimation procedure using EM algorithm has been assessed via simulation studies for both continuous and binary longitudinal data with survival time. The proposed estimates performed well in finite samples. The variance estimates based on the observed information matrix approximate the true variance well in finite samples. Simulation studies indicated that, when the true distribution of random effects is not normal, mixture distributions yield less biased estimates than no mixture and all the estimated density plots of random effects based on mixture distributions appear to have similar shapes to the true distribution. Furthermore, simulation studies also showed that the number of mixture distributions is properly selected by AIC and BIC. The proposed method was applied to data from the CHANCE study.

Our proposed method is an effort to relax the assumption that the random effects come from a normal distribution which is often made for computational reasons. This normality assumption is difficult to check because random effects are latent and never observed. Furthermore, if this assumption fails to hold, the estimates of the parameters in the generalized linear mixed model and the hazards model are biased. In this paper, the mixture of normal distributions has been shown to be a good approximation for the random effects in the simultaneous modeling when the underlying distribution of random effects is unknown. The advantage of this approach is that many continuous distributions can be well approximated by a finite normal mixture which implies that our proposed method will generally perform well. When sample size and the number of observations per subject are too large, computation may be intensive due to the integration of complete data likelihood over random effects. It will be of interest to develop a more computationally efficient approach. One possibility is to consider a penalized likelihood approach by the Laplace approximation which is currently under investigation by us.

# Chapter 5

## PENALIZED LIKELIHOOD APPROACH FOR JOINT ANALYSIS OF SURVIVAL TIME AND LONGITUDINAL OUTCOMES

### 5.1 Introduction

In biomedical or public health research, it is common that both longitudinal outcomes over time and survival endpoint are collected for the same subject along with the subject's characteristics or risk factors. Investigators are interested in finding important variables which predict both longitudinal outcomes and survival time. Among the existing approaches for longitudinal data and survival time, the selection model and the pattern mixture model have been widely used. The selection model estimating the distribution of survival time given longitudinal data was studied by numerous authors, for example, Tsiatis, Degruetola, and Wulfsohn (1995), Tsiatis and Davidian (2001), Xu and Zeger (2001a,b) and Tseng, Hsieh and Wang (2005). The pattern mixture model focuses on the trend of longitudinal outcomes conditional on survival time and was studied by Wu and Carroll (1988), Hogan and Laird (1997), Albert and Follmann (2000, 2007) and Ding and Wang (2008) among others. On the other hand, simultaneous

modeling of the longitudinal and survival data was proposed by Xu and Zeger (2001b), Zeng and Cai (2005), Elashoff, Li and Ni (2007, 2008) and Rizopoulos, Verbeke, Lesaffre and Vanrenterghem (2008). Wang and Taylor (2001), Brown and Ibrahim (2003) and Hu, Li and Li (2009) studied simultaneous modeling in the Bayesian perspective.

In the joint models, random effects are incorporated to accommodate the latent dependence between survival time and longitudinal outcomes, and often assumed to be normally distributed so that we can integrate a complete data likelihood over random effects to obtain a full likelihood. The maximum likelihood approach using an Expectation-Maximization algorithm provides the estimators which are asymptotically consistent and follows an asymptotic Gaussian process. However, the EM algorithm may be intensive on computation with large sample sizes and large numbers of longitudinal observations per subject. In the view of the cumbersome and often intractable numerical integrations required for a full likelihood, one possible alternative can be the penalized likelihood approach which gives a penalty for regarding random effects as fixed effects in the likelihood obtained by Laplace approximation. In generalized linear mixed models (GLMM), the penalized quasi-likelihood (PQL) approach is the most common estimation procedure. The PQL was proposed as an approximate Bayes procedure for some commonly occurring GLMM's by Laird (1978) and the PQL method exploited by Green (1987) for semiparametric regression analysis is available for inference in hierarchical models where the focus is on shrinkage estimation of the random effects (Robinson, 1991). Breslow and Clayton (1993) proposed to use the PQL with some modifications to a Laplace expansion for a GLMM in order to motivate standard estimating equations that may be solved by iterative application of normal theory procedures. Breslow and Lin (1995) and Lin and Brelsow (1996) derived the general expressions for the asymptotic biases in approximate estimators of regression coefficients and variance component in the GLMMs with a single source of extraneous variation

and multiple components of dispersion, respectively. The PQL also has been studied in a wide variety of GLMMs by Bartlett and Sutradhar (1999), Huber and Victoria-Feser (2004), Localio, Berlin and Ten Have (2006), Nelson and Leroux (2008), Dang, Mazumdar and Houck (2008), Jang and Lim (2009), and Masaoud and Stryhn (2010). Furthermore, the PQL is already built in SAS GLIMMIX procedure and used for the analysis of the GLMM. On the other hand, Ripatti and Palmgren (2000) proposed a penalized partial likelihood for multivariate frailty models in survival analysis. In joint modeling framework, Ye, Lin and Taylor (2008) proposed a penalized joint likelihood for a selection model and considered a continuous longitudinal process to be included as a covariate for survival time. Their penalized joint likelihood is obtained by replacing the full survival likelihood with a partial likelihood in the Laplace approximation to the full joint likelihood function, which is not equal to the actual form derived from the full joint likelihood function. On the other hand, there is no work done on the penalized likelihood approach for the simultaneous modeling of longitudinal outcomes and survival time. Furthermore, the previous study using the penalized likelihood in joint analysis (Ye, Lin and Taylor, 2008) considered continuous longitudinal data from a normal distribution.

In this paper, we propose to use a penalized likelihood to develop a more efficient estimation procedure on computation for simultaneous modeling than the EM algorithm of the maximum likelihood approach. We consider a generalized linear mixed model for longitudinal outcomes to incorporate both categorical and continuous data and a stratified Cox proportional hazards model for survival time. In this estimation procedure, all the parameters are estimated together at the same time. If the EM algorithm of maximum likelihood approach performs similarly to the penalized likelihood method on computational time, it will be better to use the full likelihood. In the meantime, if the penalized likelihood method takes less time and provides unbiased and consistent

estimates similar to those from EM algorithm, the penalized likelihood method will be preferred.

The organization of this paper is as follows. We present a simultaneous modeling for longitudinal outcomes and survival time with random effects in Section 5.2 and describe the proposed estimation procedure in Section 5.3. Numerical results from simulation studies are given in Section 5.4, and our proposed method is illustrated with the data from the Carolina Head and Neck Cancer Study (CHANCE) in Section 5.5. In Section 5.6, we discuss some further consideration.

## 5.2 Model Formulation and Notation

We use  $Y(t)$  to denote the value of a longitudinal marker process at time  $t$ . Suppose  $Y(t)$  is from a distribution belonging to exponential family in order to incorporate both continuous and categorical measurements. Let  $T$  denote survival time, and suppose that the survival time  $T$  is possibly right censored. Suppose a set of  $n$  subjects are followed over an interval  $[0, \tau]$ , where  $\tau$  is the study end time. Denote  $\mathbf{b}_i$ ,  $i = 1, \dots, n$ , as a vector of subject-specific random effects of dimension  $d_b$  and  $\mathbf{b}_i$ 's are mutually independent and identically distributed from a multivariate normal with mean zero and covariance matrix  $\Sigma_b$ .

Given the random effects  $\mathbf{b}_i$ , the observed covariates, and the observed outcome history till time  $t$ , we assume that the longitudinal outcome  $Y_i(t)$  at time  $t$  for subject  $i$  follows a distribution from the exponential family with density,

$$\exp \left\{ \frac{y_i \eta_i(t) - B(\eta_i(t))}{A(D_i(t; \phi))} + C(y_i, D_i(t; \phi)) \right\} \quad (5.1)$$

with  $\mu_i(t) = E(Y_i(t) | \mathbf{b}_i) = B'(\eta_i(t))$  and  $v_i(t) = \text{Var}(Y_i(t) | \mathbf{b}_i) = B''(\eta_i(t))A(D_i(t; \phi))$ , satisfying

$$\eta_i(t) = g(\mu_i(t)) = \mathbf{X}_i(t)\boldsymbol{\beta} + \tilde{\mathbf{X}}_i(t)\mathbf{b}_i$$

and  $v_i(t) = v(\mu_i(t))A(D_i(t; \phi))$ , where  $g(\cdot)$  and  $v(\cdot)$  are known link and variance functions respectively,  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  are the row vectors of the observed covariates for subject  $i$ , and  $\boldsymbol{\beta}$  is a column vector of coefficients for  $\mathbf{X}_i(t)$ . The random effect  $\mathbf{b}_i$  is allowed to differ for different individuals. Additionally,  $\mathbf{X}_i(t)$  and  $\tilde{\mathbf{X}}_i(t)$  can be completely different or share some components, and may include dummy variables for different strata.

Given the random effects  $\mathbf{b}_i$ , the observed covariates, and the observed survival history before time  $t$ , the conditional hazard rate function for the survival time  $T_i$  of subject  $i$  is assumed to follow a stratified multiplicative hazards model,

$$\lambda_s(t) \exp\{\tilde{\mathbf{Z}}_i(t)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(t)\boldsymbol{\gamma}\}, \quad (5.2)$$

where  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  are the row vectors of the observed covariates and may share some components,  $\boldsymbol{\psi}$  is a vector of parameters of the coefficients for random effects,  $\lambda_s(t)$  is the  $s$ -th stratum baseline hazard rate function, and  $\boldsymbol{\gamma}$  is a column vector of coefficients for  $\mathbf{Z}_i(t)$ . Note that  $\mathbf{Z}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  do not include dummy variables for strata since baseline hazard rate is stratum-specific. Here, for any vectors  $\mathbf{a}_1$  and  $\mathbf{a}_2$  of the same dimension,  $\mathbf{a}_1 \circ \mathbf{a}_2$  denotes the component-wise product. In addition,  $\tilde{\mathbf{X}}_i(t)$  and  $\tilde{\mathbf{Z}}_i(t)$  have the same dimensions as  $\mathbf{b}_i$ 's.

Under models (5.1) and (5.2), the two outcomes  $Y(t)$  and  $T$  are independent conditional on the covariates and random effect. The parameter  $\boldsymbol{\psi}$  in model (5.2) characterizes the dependence between the longitudinal outcomes and the survival time due to latent random effect:  $\boldsymbol{\psi} = \mathbf{0}$  means that the dependence between the survival time and longitudinal responses are not due to these latent variables;  $\boldsymbol{\psi} \neq \mathbf{0}$  means that such de-



pendence may be due to these latent variables. In other words,  $\boldsymbol{\psi} > \mathbf{0}$  implies that there may be some latent factors increasing both the longitudinal outcomes and the risk of survival endpoint simultaneously while  $\boldsymbol{\psi} < \mathbf{0}$  implies that some latent factors causing the increment of longitudinal outcomes may decrease the risk of survival endpoint.

Let  $n_i$  be the number of the observed longitudinal measurements for subject  $i$ , and assume that the distributions of  $n_i$  and the observation times for longitudinal measurements are independent of the parameters of interest in this joint model. The observed data from  $n$  subjects are  $(n_i, Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$ ,  $j = 1, \dots, n_i$ ,  $i = 1, \dots, n$ , and  $(V_i, \Delta_i, S_i, \{(\mathbf{Z}_i(t), \tilde{\mathbf{Z}}_i(t)) : t \leq V_i\})$ ,  $i = 1, \dots, n$ , where for subject  $i$ ,  $(Y_{ij}, \mathbf{X}_{ij}, \tilde{\mathbf{X}}_{ij})$  is the  $j$ -th observation of  $(Y_i(t), \mathbf{X}_i(t), \tilde{\mathbf{X}}_i(t))$ ,  $C_i$  is the right-censoring time and independent of  $T_i$  and  $Y_i(t)$  given the covariates and the random effects,  $V_i = \min(T_i, C_i)$ ,  $S_i$  denotes the stratum, and  $\Delta_i = I(T_i \leq C_i)$ .

Our goal is to estimate and make inferences on the parameters  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \boldsymbol{\phi}^T, \text{Vec}(\boldsymbol{\Sigma}_b)^T, \boldsymbol{\psi}^T, \boldsymbol{\gamma}^T)^T$  and the baseline cumulative hazard functions with  $S$  strata,  $\boldsymbol{\Lambda}(t) = (\Lambda_1(t), \dots, \Lambda_S(t))^T$ , where  $\Lambda_s(t) = \int_0^t \lambda_s(u) du$ ,  $s = 1, \dots, S$ . The parameters  $\boldsymbol{\beta}$  and  $\boldsymbol{\phi}$  are from the longitudinal model,  $\boldsymbol{\psi}$  and  $\boldsymbol{\gamma}$  are from the hazard model, and  $\boldsymbol{\Sigma}_b$  is associated with the random effects.  $\text{Vec}(\cdot)$  operator creates a column vector from a matrix by stacking the diagonal and upper-triangle elements of the matrix.

### 5.3 Estimation Procedure

For all  $n$  subjects, we write  $\mathbf{Y} = (\mathbf{Y}_1^T, \dots, \mathbf{Y}_n^T)^T$ ,  $\mathbf{Y}_i = (Y_{i1}, \dots, Y_{in_i})^T$ ,  $\mathbf{V} = (V_1, \dots, V_n)^T$ , and  $\mathbf{b} = (\mathbf{b}_1^T, \dots, \mathbf{b}_n^T)^T$ . Then, the likelihood function of the complete data  $(\mathbf{Y}, \mathbf{V}, \mathbf{b})$  has the form,

$$\begin{aligned}
L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}) &= \prod_{s=1}^S \prod_{i=1}^n [f(\mathbf{Y}_i, V_i | \mathbf{b}_i) f(\mathbf{b}_i)]^{I(S_i=s)} = \prod_{i=1}^n f(\mathbf{Y}_i | \mathbf{b}_i) \left( \prod_{s=1}^S [f(V_i | \mathbf{b}_i)]^{I(S_i=s)} \right) f(\mathbf{b}_i) \\
&= \prod_{i=1}^n \exp \left\{ \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij} \boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij} \mathbf{b}_i) - B(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \right\} \\
&\quad \times \left( \prod_{s=1}^S \left[ \lambda_s(V_i)^{\Delta_i} \exp \left\{ \Delta_i [\widetilde{\mathbf{W}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{W}_i(V_i)\boldsymbol{\gamma}] \right. \right. \right. \\
&\quad \left. \left. \left. - \int_0^{V_i} \exp \{ \widetilde{\mathbf{W}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{W}_i(u)\boldsymbol{\gamma} \} d\Lambda_s(u) \right\} \right]^{I(S_i=s)} \right) \\
&\quad \times (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right\},
\end{aligned}$$

and the full likelihood function of the observed data  $(\mathbf{Y}, \mathbf{V})$  for the parameter  $(\boldsymbol{\theta}, \boldsymbol{\Lambda})$  is expressed as

$$L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = \int_{\mathbf{b}} L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b}) d\mathbf{b}. \quad (5.3)$$

The primary difficulty in implementing this full likelihood inference lies in the integrations needed to evaluate the complete data likelihood  $L_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  and its partial derivatives.

In the EM algorithm of maximum likelihood approach, the random effect  $\mathbf{b}_i$  is considered as missing data for  $i = 1, \dots, n$ . Thus, the M-step solves the conditional score equations from complete data log-likelihood given observations, where the conditional expectation is evaluated in the E-step. The procedure involves iterating between the two steps until convergence is achieved. In the E-step calculating the conditional expectations of some known functions of  $\mathbf{b}_i$  needed in the next M-step, a numerical approximation method such as the Gauss-Hermite Quadrature is required for the integration with the posterior probability of random effects. When sample size ( $n$ ), the

number of observations per subject ( $n_i$ ), and the number of parameters to be estimated are large, the task involving the integration in the E-step is very challenging. Therefore, to make the simultaneous modeling more practical, we build the algorithm which relieves the computational burden.

Our proposed estimation method is to calculate the maximum penalized likelihood estimates for  $(\boldsymbol{\theta}, \boldsymbol{\Lambda}(t))$  over a set in which  $\boldsymbol{\theta}$  is in a bounded set and  $\Lambda_s(t)$  of  $\boldsymbol{\Lambda}(t)$  belongs to a space consisting of all the increasing functions with  $\Lambda_s(0) = 0$ ,  $s = 1, \dots, S$ . We let each  $\Lambda_s(t)$  of  $\boldsymbol{\Lambda}(t)$ ,  $s = 1, \dots, S$ , be an increasing and right-continuous step function with jumps only at the observed failure times belonging to stratum  $s$ . The penalized likelihood is obtained by Laplace approximation, and the proposed approach is expected to be less intensive in computation in the sense that it imposes the penalty for considering the random effect as the fixed effect in the likelihood and therefore no calculation for integrating the likelihood over random effects is needed.

### 5.3.1 Laplace approximation

The full likelihood (5.3) can be written as

$$L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = (2\pi)^{-nd_b/2} |\boldsymbol{\Sigma}_b|^{-n/2} \int_{\mathbf{b}} \exp \left\{ \sum_{i=1}^n \left[ l_{i|b}(\boldsymbol{\theta}, \Lambda_s) - \frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right] \right\} d\mathbf{b}, \quad (5.4)$$

where the logarithm of the conditional joint density given an unobserved random effect  $\mathbf{b}_i$  is

$$\begin{aligned} l_{i|b}(\boldsymbol{\theta}, \Lambda_s) = & \sum_{j=1}^{n_i} \left[ \frac{Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i) - B_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} + C(Y_{ij}; D_i(t_j; \phi)) \right] \\ & + \sum_{s=1}^S I(S_i = s) \left[ \Delta_i \log(\lambda_s(V_i)) + \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \\ & \left. - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} d\Lambda_s(u) \right]. \quad (5.5) \end{aligned}$$

Then, we have the following form of the full log-likelihood,

$$l_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) = \sum_{i=1}^n \left[ -\frac{d_b}{2} \log(2\pi) - \frac{1}{2} |\boldsymbol{\Sigma}_b| + l_{i|b}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) - \frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right]. \quad (5.6)$$

In (5.4), define

$$-\boldsymbol{\kappa}(\mathbf{b}) = \sum_{i=1}^n \left[ l_{i|b}(\boldsymbol{\theta}, \boldsymbol{\Lambda}_s) - \frac{1}{2} \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \mathbf{b}_i \right] = \sum_{i=1}^n \left[ -\boldsymbol{\kappa}_i(\mathbf{b}_i) \right] \quad (5.7)$$

and apply Laplace's approximation as following,

$$-\boldsymbol{\kappa}_i(\mathbf{b}_i) \approx -\boldsymbol{\kappa}_i(\tilde{\mathbf{b}}_i) - \frac{1}{2} (\mathbf{b}_i - \tilde{\mathbf{b}}_i)^T \boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i) (\mathbf{b}_i - \tilde{\mathbf{b}}_i),$$

where  $\boldsymbol{\kappa}'$  and  $\boldsymbol{\kappa}''$  denote the  $d_b$  vector and  $d_b \times d_b$  dimensional matrix of first- and second-order partial derivatives of  $\boldsymbol{\kappa}$  with respect to  $\mathbf{b}$  and  $\tilde{\mathbf{b}}$  denotes the solution to  $\boldsymbol{\kappa}'(\mathbf{b}) = 0$  that minimizes  $\boldsymbol{\kappa}(\mathbf{b})$ . Then, the full likelihood function (5.4) can be approximated as followings,

$$\begin{aligned} & L_f(\boldsymbol{\theta}, \boldsymbol{\Lambda}; \mathbf{Y}, \mathbf{V}) \\ &= (2\pi)^{-nd_b/2} |\boldsymbol{\Sigma}_b|^{-n/2} \int_{\mathbf{b}} \exp \left\{ -\boldsymbol{\kappa}(\mathbf{b}) \right\} d\mathbf{b} = \prod_{i=1}^n \left[ (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \int_{\mathbf{b}} \exp \left\{ -\boldsymbol{\kappa}_i(\mathbf{b}_i) \right\} d\mathbf{b} \right] \\ &\approx \prod_{i=1}^n \left[ (2\pi)^{-d_b/2} |\boldsymbol{\Sigma}_b|^{-1/2} \int_{\mathbf{b}} \exp \left\{ -\boldsymbol{\kappa}_i(\tilde{\mathbf{b}}_i) - \frac{1}{2} (\mathbf{b}_i - \tilde{\mathbf{b}}_i)^T \boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i) (\mathbf{b}_i - \tilde{\mathbf{b}}_i) \right\} d\mathbf{b} \right] \\ &= \prod_{i=1}^n \left[ |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\boldsymbol{\kappa}_i(\tilde{\mathbf{b}}_i) \right\} |\boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i)|^{-1/2} \right. \\ &\quad \left. (2\pi)^{-d_b/2} |\boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i)|^{1/2} \int_{\mathbf{b}} \exp \left\{ -\frac{1}{2} (\mathbf{b}_i - \tilde{\mathbf{b}}_i)^T \boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i) (\mathbf{b}_i - \tilde{\mathbf{b}}_i) \right\} d\mathbf{b} \right] \\ &= \prod_{i=1}^n \left[ |\boldsymbol{\Sigma}_b|^{-1/2} \exp \left\{ -\boldsymbol{\kappa}_i(\tilde{\mathbf{b}}_i) \right\} |\boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i)|^{-1/2} \right] \\ &= |\boldsymbol{\Sigma}_b|^{-n/2} \exp \left\{ \sum_{i=1}^n \left[ -\boldsymbol{\kappa}_i(\tilde{\mathbf{b}}_i) - \frac{1}{2} \log |\boldsymbol{\kappa}_i''(\tilde{\mathbf{b}}_i)| \right] \right\}. \end{aligned} \quad (5.8)$$

In the above (5.8), ignoring the multiplicative constant, the approximation yields

$$-\kappa(\mathbf{b}) \approx -\frac{1}{2} \log |\kappa''(\tilde{\mathbf{b}})| - \kappa(\tilde{\mathbf{b}}).$$

Note that, from (5.7),

$$\begin{aligned}\kappa_i(\tilde{\mathbf{b}}_i) &= -\tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s) + \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i, \\ \kappa'_i(\tilde{\mathbf{b}}_i) &= -\tilde{l}'_{i|b}(\boldsymbol{\theta}, \Lambda_s) + \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i, \\ \kappa''_i(\tilde{\mathbf{b}}_i) &= -\tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s) + \boldsymbol{\Sigma}_b^{-1},\end{aligned}\tag{5.9}$$

where  $\tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s)$  is (5.5) evaluated at  $\tilde{\mathbf{b}}_i$ , and  $\tilde{l}'_{i|b}(\boldsymbol{\theta}, \Lambda_s)$  and  $\tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)$  are the first and second derivatives of (5.5) with respect to  $\mathbf{b}_i$  evaluated at  $\tilde{\mathbf{b}}_i$ . Then, the first order Laplace approximation to the full likelihood becomes

$$\begin{aligned}(5.8) &= \exp \left\{ \sum_{i=1}^n \left[ -\frac{1}{2} \log |\boldsymbol{\Sigma}_b| - \kappa_i(\tilde{\mathbf{b}}_i) - \frac{1}{2} \log |-\tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s) + \boldsymbol{\Sigma}_b^{-1}| \right] \right\} \\ &= \exp \left\{ \sum_{i=1}^n \left[ -\frac{1}{2} \log |\mathbf{I}_{d_b} - \boldsymbol{\Sigma}_b \tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)| + \tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s) - \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i \right] \right\} \\ &= \exp \left\{ \sum_{i=1}^n \tilde{l}_{fi}(\boldsymbol{\theta}, \Lambda_s) \right\} = \exp \left\{ \tilde{l}_f(\boldsymbol{\theta}, \Lambda_s) \right\},\end{aligned}$$

where  $\tilde{l}_f(\boldsymbol{\theta}, \Lambda_s)$  is the first order Laplace approximation to the full log-likelihood function  $l_f(\boldsymbol{\theta}, \Lambda_s)$  of (5.6). That is,

$$\tilde{l}_f(\boldsymbol{\theta}, \Lambda) = \sum_{i=1}^n \left[ -\frac{1}{2} \log |\mathbf{I}_{d_b} - \boldsymbol{\Sigma}_b \tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)| + \left( \tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s) - \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i \right) \right], \tag{5.10}$$

where  $\tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s) = l''_{i|b}(\boldsymbol{\theta}, \Lambda_s; \tilde{\mathbf{b}}_i)$  and  $\tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s) = l_{i|b}(\boldsymbol{\theta}, \Lambda_s; \tilde{\mathbf{b}}_i)$ , and the first and second derivatives of  $l_{i|b}(\boldsymbol{\theta}, \Lambda_s; \tilde{\mathbf{b}}_i)$  in (5.5) with respect to  $\mathbf{b}_i$  are

$$l'_{i|b}(\boldsymbol{\theta}, \Lambda_s) = \sum_{j=1}^{n_i} \left[ \frac{Y_{ij} \tilde{\mathbf{X}}_{ij}}{A(D_i(t_j; \phi))} - \frac{B'_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} \right] + \sum_{s=1}^S I(S_i=s) \left[ \Delta_i(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\}(\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right]$$

and

$$l''_{i|b}(\boldsymbol{\theta}, \Lambda_s) = - \sum_{j=1}^{n_i} \frac{B''_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} - \sum_{s=1}^S I(S_i=s) \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\}(\tilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi})^{\otimes 2} d\Lambda_s(u), \quad (5.11)$$

where  $B'_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)$  and  $B''_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)$  are the first and second derivatives of  $B_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)$  with respect to  $\mathbf{b}_i$ .

### 5.3.2 Penalized likelihood

Now, we further approximate (5.11). The first term of (5.11) can be expressed as  $\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i$ , where  $\mathbf{W}_i$  is the  $n_i \times n_i$  diagonal matrix with  $w_{ij} = A(D_i(t_j; \phi)) [g'(\mu_{ij}^b)]^{-1}$ ,  $g(\cdot)$  is a canonical link function,  $\mu_{ij}^b = E(Y_{ij} | \mathbf{b}_i)$ ,  $g'(\mu_{ij}^b)$  is the derivative of  $g(\mu_{ij}^b)$  with respect to  $\mu_{ij}^b$ , and  $\tilde{\mathbf{X}}_i = (\tilde{\mathbf{X}}_{i1}^T, \dots, \tilde{\mathbf{X}}_{in_i}^T)^T$ . Generalized linear model (GLM) iterative weights  $\mathbf{W}_i$  (i.e.  $w_{ij}$ ) vary slowly or not at all at the function of the mean, and hance, by taking an expectation,

$$E \left[ - \sum_{j=1}^{n_i} \frac{B''_{ij}(\boldsymbol{\beta}; \mathbf{b}_i)}{A(D_i(t_j; \phi))} \right] = E [\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i] \quad (5.12)$$

becomes a constant. In the second term of (5.11), by taking an expectation, we have

$$E \left[ \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\}(\tilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right]$$

$$\begin{aligned}
&= \text{E} \left[ \int (\tilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) I(V_i \geq u) \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} \lambda_s(u) du \right] \\
&= \text{E} \left[ \int (\tilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) I(T_i \geq u) I(C_i \geq u) \frac{1}{S_{sT}(u)} dF_{sT}(u) \right] \\
&= \text{E} \left[ (\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \Delta_i \right]
\end{aligned}$$

becomes a constant. Thus, the expectation of second term in (5.11)

$$\sum_{s=1}^S I(S_i = s) \text{E} \left[ -(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \Delta_i \right] \quad (5.13)$$

also becomes a constant including  $\boldsymbol{\psi}$ . Since the expected value of (5.11),  $\text{E}[l''_{i|b}(\boldsymbol{\theta}, \Lambda_s)]$ , is the sum of (5.12) and (5.13),  $l''_{i|b}(\boldsymbol{\theta}, \Lambda_s)$  in (5.11) also becomes asymptotically a constant. Therefore, in (5.10),  $|\mathbf{I}_{db} - \boldsymbol{\Sigma}_b \tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)|$  is asymptotically a constant including  $\boldsymbol{\Sigma}_b$  and  $\boldsymbol{\psi}$ . Then, we derive the penalized log-likelihood as following,

$$\begin{aligned}
&l_P(\boldsymbol{\theta}, \Lambda) \\
&= \sum_{i=1}^n \left[ -\frac{1}{2} \log \left| \mathbf{I}_{db} - \boldsymbol{\Sigma}_b \left( \text{E}[\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i] - \sum_{s=1}^S I(S_i = s) \text{E}[(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \Delta_i] \right) \right| \right. \\
&\quad \left. + \left( \tilde{l}_{i|b}(\boldsymbol{\theta}, \Lambda_s) - \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i \right) \right]. \quad (5.14)
\end{aligned}$$

Since  $|\mathbf{I}_{db} - \boldsymbol{\Sigma}_b \tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)|$  in (5.10) which corresponds to the first term in (5.14) is asymptotically a constant including  $\boldsymbol{\Sigma}_b$  and  $\boldsymbol{\psi}$ , it only contributes to the estimating equations of  $\boldsymbol{\Sigma}_b$  and  $\boldsymbol{\psi}$ , and we ignore the term to obtain the estimating equations of  $\boldsymbol{\beta}$ ,  $\boldsymbol{\phi}$  and  $\boldsymbol{\gamma}$ .  $-\frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i$  is the penalty term for regarding the random effects as fixed effects by replacing  $\mathbf{b}$  with  $\tilde{\mathbf{b}}$  in the likelihood. We choose  $\boldsymbol{\theta}$  to maximize the penalized likelihood  $l_p(\boldsymbol{\theta}, \Lambda)$  in (5.14). That is,  $(\hat{\boldsymbol{\theta}}, \tilde{\mathbf{b}})$  jointly maximize the equation (5.14). The score equations for  $(\boldsymbol{\theta}, \mathbf{b})$  are obtained by differentiating (5.14) with respect to  $\boldsymbol{\theta}$  and  $\mathbf{b}$ , respectively.

### 5.3.3 Implementation

We conduct the Newton-Rapshon method for estimating equations to obtain  $\tilde{\mathbf{b}}$  and  $\widehat{\boldsymbol{\theta}}$ . The procedure involves iterating between the following two steps until convergence is achieved: at the  $k$ -th iteration,

Step1 : Conduct one-step Newton-Rapshon iteration to obtain the solution  $\tilde{\mathbf{b}}$  of  $\boldsymbol{\kappa}'(\mathbf{b}) = 0$ . The  $(k+1)$ -th estimate is  $\tilde{\mathbf{b}}^{(k+1)} = \tilde{\mathbf{b}}^{(k)} - [\boldsymbol{\kappa}''(\tilde{\mathbf{b}}^{(k)})]^{-1}[\boldsymbol{\kappa}'(\tilde{\mathbf{b}}^{(k)})]^T$ , where  $\tilde{\mathbf{b}}^{(k)} = \tilde{\mathbf{b}}^{(k)}(\widehat{\boldsymbol{\theta}}^{(k-1)})$ ,  $\boldsymbol{\kappa}'(\mathbf{b}) = (\boldsymbol{\kappa}'_1(\mathbf{b}_1)^T, \dots, \boldsymbol{\kappa}'_n(\mathbf{b}_n)^T)^T$  and  $\boldsymbol{\kappa}''(\mathbf{b}) = (\boldsymbol{\kappa}''_1(\mathbf{b}_1)^T, \dots, \boldsymbol{\kappa}''_n(\mathbf{b}_n)^T)^T$ , and the functions  $\boldsymbol{\kappa}'_i(\mathbf{b}_i)$  and  $\boldsymbol{\kappa}''_i(\mathbf{b}_i)$ ,  $i = 1, \dots, n$ , are given in (5.9).

Step2 : By one-step Newton-Rapshon iteration, the  $(k+1)$ -th estimate is calculated as  $\widehat{\boldsymbol{\theta}}^{(k+1)} = \widehat{\boldsymbol{\theta}}^{(k)} - [S'_P(\widehat{\boldsymbol{\theta}}^{(k)})]^{-1}[S_P(\widehat{\boldsymbol{\theta}}^{(k)})]^T$ , where  $S_P(\boldsymbol{\theta})$  is the score equation for  $\boldsymbol{\theta}$  from the penalized log-likelihood and  $S'_P(\boldsymbol{\theta})$  is the first derivative of  $S_P(\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ . With  $(\widehat{\boldsymbol{\theta}}^{(k+1)}, \tilde{\mathbf{b}}^{(k+1)})$ , the  $(k+1)$ -th Breslow-type estimate of the baseline cumulative hazard for the  $s$ -th stratum is obtained as an empirical function which has jumps only at the observed failure time,

$$\begin{aligned} \Lambda_s^{(k+1)}(t) &= \Lambda_s^{(k+1)}(t; \widehat{\boldsymbol{\theta}}^{(k+1)}, \tilde{\mathbf{b}}^{(k+1)}) \\ &= \sum_{i: V_i \leq t} \frac{\Delta_i I(S_i = s)}{\sum_{l: V_l \geq V_i} \exp\{\tilde{\mathbf{Z}}_l(V_i)(\widehat{\boldsymbol{\psi}}^{(k+1)} \circ \tilde{\mathbf{b}}_l^{(k+1)}) + \mathbf{Z}_l(V_i)\widehat{\boldsymbol{\gamma}}^{(k+1)}\}} I(S_l = s)}. \end{aligned} \quad (5.15)$$

For variance estimation of  $(\widehat{\boldsymbol{\theta}}, \widehat{\Lambda}(t))$ , we adopt the observed information matrix via Louis (1982) formula and conduct the Expectation step used in the maximum likelihood approach with the estimates by the penalized likelihood method. For the numerical calculation of the observed information matrix, we consider  $\Lambda_s\{V_i\}$ , the jump size of  $\Lambda_s(t)$  at  $V_i$  belonging to stratum  $s$  for which  $\Delta_i = 1$ , instead of  $\lambda_s(V_i)$ . That is,  $\boldsymbol{\Lambda}\{\cdot\} = (\boldsymbol{\Lambda}_1^T\{\cdot\}, \dots, \boldsymbol{\Lambda}_S^T\{\cdot\})^T$  with  $\boldsymbol{\Lambda}_s\{\cdot\} = (\Lambda\{T_{s1}\}, \dots, \Lambda\{T_{sm_s}\})^T$  for  $m_s$  failure times among  $n_s$  subjects ( $0 \leq m_s \leq n_s$ ) of the  $s$ -th stratum,  $s = 1, \dots, S$ . Then, by the Louis (1982) formula,



$$\begin{aligned}
I(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}) &= E_{b|Y,V}[B_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})|\mathbf{Y}, \mathbf{V}] \\
&- E_{b|Y,V}[U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})U_c^T(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})|\mathbf{Y}, \mathbf{V}] \\
&+ E_{b|Y,V}[U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})] E_{b|Y,V}[U_c^T(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})],
\end{aligned}$$

where  $U_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  and  $B_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$  are the first derivative vector and the negative of the second derivative matrix for the complete data log-likelihood  $l_c(\boldsymbol{\theta}, \boldsymbol{\Lambda}\{\cdot\}; \mathbf{Y}, \mathbf{V}, \mathbf{b})$ , respectively. For subject  $i$  with  $S_i = s$ , given observations and the penalized likelihood estimate  $(\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s)$ , we calculate the following conditional expectation of a known function  $q(\mathbf{b}_i)$  needed in the observed information matrix,

$$\begin{aligned}
E[q(\mathbf{b}_i)|\widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s] &= \frac{\int \mathbf{b}_i q(\mathbf{b}_i) f(\mathbf{Y}_i, V_i | \mathbf{b}_i, \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s) f(\mathbf{b}_i | \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s) d\mathbf{b}_i}{\int \mathbf{b}_i f(\mathbf{Y}_i, V_i | \mathbf{b}_i, \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s) f(\mathbf{b}_i | \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s) d\mathbf{b}_i} \\
&= \frac{\int \mathbf{z}_G q(R(\mathbf{z}_G)) K(\mathbf{z}_G) \exp\{-\mathbf{z}_G^T \mathbf{z}_G\} d\mathbf{z}_G}{\int \mathbf{z}_G K(\mathbf{z}_G) \exp\{-\mathbf{z}_G^T \mathbf{z}_G\} d\mathbf{z}_G}, \tag{5.16}
\end{aligned}$$

where  $\mathbf{z}_G$  follows a multivariate Gaussian distribution with mean zero,  $\mathbf{z}_G = R^{-1}(\mathbf{b}_i)$ ,  $K(\mathbf{z}_G) = \exp\{\mathbf{z}_G^T \mathbf{z}_G\} f(\mathbf{Y}_i, V_i | R(\mathbf{z}_G), \boldsymbol{\theta}^{(k)}, \boldsymbol{\Lambda}_s^{(k)}) f(R(\mathbf{z}_G) | \widehat{\boldsymbol{\theta}}, \widehat{\boldsymbol{\Lambda}}_s)$ , and Gauss-Hermite Quadrature numerical approximation is used for the calculation of integration. Note that, in (5.16), the functions of  $R(\cdot)$  and  $K(\cdot)$  have different expressions for different longitudinal distributions.

The proposed penalized likelihood approach for simultaneous modeling can be applied to all generalized linear mixed models of longitudinal outcomes. Next, we provide the expressions of the penalized log-likelihood and relevant equations for continuous and binary longitudinal outcomes with survival time.

*Ex 1. Continuous longitudinal data with Normal distribution and survival time*

Continuous longitudinal outcomes following a normal distribution has  $A(D_i(t_j; \phi)) =$

$\sigma_y^2$ ,  $B_{ij}(\beta; \mathbf{b}_i) = (\mathbf{X}_{ij}\beta + \tilde{\mathbf{X}}_{ij}\mathbf{b}_i)^2/2$ , and  $C(Y_{ij}; D_i(t_j; \phi)) = -(y^2/\sigma_y^2 + \log(\sigma_y^2) + \log(2\pi))/2$  in (5.1), where  $\sigma_y^2$  is the variance of longitudinal outcomes given  $\mathbf{b}_i$ . Then, the  $\kappa'_i(\mathbf{b}_i)$  and  $\kappa''_i(\mathbf{b}_i)$ ,  $i = 1, \dots, n$ , used in Step 1 are

$$\begin{aligned} \kappa'_i(\mathbf{b}_i) &= - \left[ \sum_{j=1}^{n_i} \frac{1}{\sigma_y^2} (Y_{ij} - \mathbf{X}_{ij}\beta - \tilde{\mathbf{X}}_{ij}\mathbf{b}_i) \tilde{\mathbf{X}}_{ij} + \sum_{s=1}^S I(S_i=s) \left( \Delta_i(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \right. \right. \\ &\quad \left. \left. - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} (\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right) \right] + \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \end{aligned}$$

and

$$\begin{aligned} \kappa''_i(\mathbf{b}_i) &= - \left[ \sum_{j=1}^{n_i} \left( -\frac{1}{\sigma_y^2} \right) \tilde{\mathbf{X}}_{ij}^T \tilde{\mathbf{X}}_{ij} \right. \\ &\quad \left. + \sum_{s=1}^S I(S_i=s) \left( - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} (\tilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi}) (\tilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right) \right] \\ &\quad + \boldsymbol{\Sigma}_b^{-1}. \end{aligned}$$

In Step2, the penalized log-likelihood (5.14) has the following form for continuous longitudinal outcomes from a normal distribution and survival time,

$$\begin{aligned} l_P(\boldsymbol{\theta}, \boldsymbol{\Lambda}) &= \sum_{i=1}^n \left[ -\frac{1}{2} \log \left| \mathbf{I}_{db} - \boldsymbol{\Sigma}_b \left( \mathbb{E}[\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i] - \sum_{s=1}^S I(S_i=s) \mathbb{E}[(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \Delta_i] \right) \right| \right. \\ &\quad - \sum_{j=1}^{n_i} \frac{1}{2\sigma_y^2} (Y_{ij} - \mathbf{X}_{ij}\beta - \tilde{\mathbf{X}}_{ij}\tilde{\mathbf{b}}_i)^2 - \frac{n_i}{2} \log(2\pi\sigma_y^2) \\ &\quad + \sum_{s=1}^S I(S_i=s) \left[ \Delta_i \log(\lambda_s(V_i)) + \Delta_i[\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \tilde{\mathbf{b}}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \\ &\quad \left. \left. - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \tilde{\mathbf{b}}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} d\Lambda_s(u) \right] - \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i \right]. \quad (5.17) \end{aligned}$$

For the expected values in (5.17), we use the original expressions of  $\mathbf{W}$  and  $\Delta_i$ , evaluated at the parameter estimates at the previous iteration and the random effects estimates from Step 1 at the current iteration, as  $\widehat{\mathbf{W}}$  and  $\widehat{\Delta}_i$ , which are

$$-\sum_{j=1}^{n_i} \frac{1}{\widehat{\sigma}_y^2} \widetilde{\mathbf{X}}_{ij}^T \widetilde{\mathbf{X}}_{ij} \quad \text{and} \quad \exp\{\widetilde{\mathbf{Z}}_i(V_i)(\widehat{\boldsymbol{\psi}} \circ \widetilde{\mathbf{b}}_i) + \mathbf{Z}_i(V_i)\widehat{\boldsymbol{\gamma}}\}\widehat{\Lambda}_s(V_i),$$

respectively. On the other hand, the observed variance of longitudinal outcomes and the observed event for each subject also may be used as  $\widehat{\sigma}_y^2$  and  $\widehat{\Delta}_i$ .  $S_P(\boldsymbol{\theta})$  is obtained by differentiating (5.17) with respect to  $\boldsymbol{\theta}$ , and  $S'_P(\boldsymbol{\theta})$  is the derivative of  $S_P(\boldsymbol{\theta})$  with respect to  $\boldsymbol{\theta}$ . The Breslow-type estimator of the baseline cumulative hazard for the  $s$ -th stratum has the same expression given in (5.15) for all different longitudinal distributions.

*Ex 2. Binary longitudinal data and survival time*

Logistic distribution has  $A(D_i(t_j; \phi)) = 1$ ,  $B_{ij}(\boldsymbol{\beta}; \mathbf{b}_i) = \log(1 + \exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_i\})$ , and  $C(Y_{ij}; D_i(t_j; \phi)) = 0$  in (5.1). Thus, the  $\boldsymbol{\kappa}'_i(\mathbf{b}_i)$  and  $\boldsymbol{\kappa}''_i(\mathbf{b}_i)$ ,  $i = 1, \dots, n$ , in Step 1 are

$$\begin{aligned} & \boldsymbol{\kappa}'_i(\mathbf{b}_i) \\ &= -\left[ \sum_{j=1}^{n_i} \left( Y_{ij} - \frac{\exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_i\}}{1 + \exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_i\}} \right) \widetilde{\mathbf{X}}_{ij} + \sum_{s=1}^S I(S_i = s) \left( \Delta_i(\widetilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \right. \right. \\ & \quad \left. \left. - \int_0^{V_i} \exp\{\widetilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\}(\widetilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right) \right] + \mathbf{b}_i^T \boldsymbol{\Sigma}_b^{-1} \end{aligned}$$

and

$$\begin{aligned} & \boldsymbol{\kappa}''_i(\mathbf{b}_i) \\ &= -\left[ \sum_{j=1}^{n_i} \left( -\frac{\exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_i\}}{(1 + \exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \widetilde{\mathbf{X}}_{ij}\mathbf{b}_i\})^2} \right) \widetilde{\mathbf{X}}_{ij}^T \widetilde{\mathbf{X}}_{ij} \right. \\ & \quad \left. + \sum_{s=1}^S I(S_i = s) \left( -\int_0^{V_i} \exp\{\widetilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \mathbf{b}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\}(\widetilde{\mathbf{Z}}_i^T(u) \circ \boldsymbol{\psi})(\widetilde{\mathbf{Z}}_i(u) \circ \boldsymbol{\psi}^T) d\Lambda_s(u) \right) \right] \\ & \quad + \boldsymbol{\Sigma}_b^{-1}. \end{aligned}$$

In Step 2, the penalized log-likelihood (5.14) has the following form for binary longitudinal outcomes and survival time,

$$\begin{aligned}
& l_P(\boldsymbol{\theta}, \boldsymbol{\Lambda}) \\
&= \sum_{i=1}^n \left[ -\frac{1}{2} \log \left| \mathbf{I}_{db} - \boldsymbol{\Sigma}_b \left( \mathbb{E} [\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i] - \sum_{s=1}^S I(S_i=s) \mathbb{E} [(\tilde{\mathbf{Z}}_i^T(V_i) \circ \boldsymbol{\psi})(\tilde{\mathbf{Z}}_i(V_i) \circ \boldsymbol{\psi}^T) \Delta_i] \right) \right| \right. \\
&\quad + \sum_{j=1}^{n_i} [Y_{ij}(\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\tilde{\mathbf{b}}_i) - \log(1 + \exp\{\mathbf{X}_{ij}\boldsymbol{\beta} + \tilde{\mathbf{X}}_{ij}\tilde{\mathbf{b}}_i\})] \\
&\quad + \sum_{s=1}^S I(S_i=s) \left[ \Delta_i \log(\lambda_s(V_i)) + \Delta_i [\tilde{\mathbf{Z}}_i(V_i)(\boldsymbol{\psi} \circ \tilde{\mathbf{b}}_i) + \mathbf{Z}_i(V_i)\boldsymbol{\gamma}] \right. \\
&\quad \left. \left. - \int_0^{V_i} \exp\{\tilde{\mathbf{Z}}_i(u)(\boldsymbol{\psi} \circ \tilde{\mathbf{b}}_i) + \mathbf{Z}_i(u)\boldsymbol{\gamma}\} d\Lambda_s(u) \right] - \frac{1}{2} \tilde{\mathbf{b}}_i^T \boldsymbol{\Sigma}_b^{-1} \tilde{\mathbf{b}}_i \right]. \quad (5.18)
\end{aligned}$$

For the expected values in (5.18), we use the original expressions of  $\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i$  and  $\Delta_i$ , evaluated at the parameter estimates at the previous iteration and the random effects estimates from Step 1 at the current iteration, as  $\tilde{\mathbf{X}}_i^T \widehat{\mathbf{W}} \tilde{\mathbf{X}}_i$  and  $\widehat{\Delta}_i$ , which are

$$-\sum_{j=1}^{n_i} \frac{\exp\{\mathbf{X}_{ij}\widehat{\boldsymbol{\beta}} + \tilde{\mathbf{X}}_{ij}\tilde{\mathbf{b}}_i\}}{(1 + \exp\{\mathbf{X}_{ij}\widehat{\boldsymbol{\beta}} + \tilde{\mathbf{X}}_{ij}\tilde{\mathbf{b}}_i\})^2} \tilde{\mathbf{X}}_{ij}^T \tilde{\mathbf{X}}_{ij} \quad \text{and} \quad \exp\{\tilde{\mathbf{Z}}_i(V_i)(\widehat{\boldsymbol{\psi}} \circ \tilde{\mathbf{b}}_i) + \mathbf{Z}_i(V_i)\widehat{\boldsymbol{\gamma}}\} \widehat{\Lambda}_s(V_i),$$

respectively. On the other hand, since the original expression of  $\tilde{\mathbf{X}}_i^T \mathbf{W}_i \tilde{\mathbf{X}}_i$  is same as  $-n_i \text{Var}(Y_{ij}|\mathbf{b}_i)$ , the observed variance of longitudinal outcomes of each subject may be used as  $\widehat{\text{Var}}(Y_{ij}|\mathbf{b}_i)$ . Likewise, the individual observed event also may be used as  $\widehat{\Delta}_i$ .  $S_P(\boldsymbol{\theta})$  and  $S'_P(\boldsymbol{\theta})$  are the first and second derivatives of (5.18) with respect to  $\boldsymbol{\theta}$ .

## 5.4 Simulation Studies

In this section, through simulation studies, we compare numerical performances on the computing time, bias, and mean squared error (MSE) of the penalized likelihood method and the EM algorithm used in maximum likelihood estimation for the simultaneous modeling of binary longitudinal outcomes and survival time with a random

intercept.

We assume that  $Y_{ij}$  is a binary outcome following

$$P(Y_{ij} = y_{ij}|b_i) = \exp \left\{ y_{ij}\eta_{ij} - \log(1 + \exp\{\eta_{ij}\}) \right\}, \quad y_{ij} = 0, 1,$$

with  $\eta_{ij} = \mathbf{X}_{ij}\boldsymbol{\beta} + b_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \beta_3 X_{3ij} + b_i$  for  $j = 1, \dots, n_i$ , and

$$h(t|b_i) = \lambda(t) \exp\{\psi b_i + \mathbf{Z}_i(t)\boldsymbol{\gamma}\} = \lambda(t) \exp\{\psi b_i + \gamma_1 Z_{1i} + \gamma_2 Z_{2i}\},$$

where  $b_i \sim N(0, \sigma_b^2)$ ,  $X_{1i} \equiv Z_{1i}$  are simulated from a Bernoulli distribution with success probability being 0.5, and  $X_{2i} \equiv Z_{2i}$  are simulated from the uniform distribution between 0 and 1. For the time at which longitudinal data are observed, we consider 4 different units of 0.3, 0.1, 0.05 and 0.03. The longitudinal data are generated for every unit of time, and thus  $X_{3ij}$ , the time at measurement, has the value of every unit ranging over 0 through 2.4. We consider  $\psi = -0.1$  indicating negative dependency between longitudinal process and survival time model. The parameters in the two models are chosen as  $\beta_0 = -1$ ,  $\beta_1 = 1$ ,  $\beta_2 = -0.5$ ,  $\beta_3 = -0.2$ ,  $\sigma_b^2 = 0.5$ ,  $\psi = -0.1$ ,  $\gamma_1 = -0.1$ ,  $\gamma_2 = 0.1$ , and  $\lambda(t) = 1$ . Censoring time is generated from the uniform distribution between 0.4 and 2.4, and the censoring proportion is around 25~35%. We consider different sample sizes ( $n$ ) of 200 and 400 with 1000 replications and different average numbers of longitudinal observations per subject ( $n_i$ ) which are 4, 8, 15 and 25. For the estimated baseline cumulative hazard function, we consider three fixed time points of 0.9, 1.4, and 1.9. Table 5.1 and Table 5.2 show the simulation results of maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) for  $\boldsymbol{\theta} = (\boldsymbol{\beta}^T, \sigma_b^2, \psi, \boldsymbol{\gamma}^T)^T$  and baseline cumulative hazards at the given three time points in the simultaneous modeling of binary longitudinal outcomes and survival time with sample sizes of 200 and 400, respectively. In Table 5.1 and Table 5.2, “ $n_i$ ” is the average number of lon-

Table 5.1: Summary of simulation results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) in the simultaneous modeling of binary longitudinal outcomes and survival time (n=200).

$n_i$	Par.	True	MLE						MPLE					
			Est.	Bias	SSD	ESE	MSE	CP	Est.	Bias	SSD	ESE	MSE	CP
4	$\beta_0$	-1.0	-1.013	-.013	.247	.243	.061	.948	-.932	.068	.226	.228	.056	.948
	$\beta_1$	1.0	1.004	.004	.203	.205	.041	.953	.928	-.072	.186	.190	.040	.933
	$\beta_2$	-.5	-.479	-.021	.358	.349	.128	.944	-.444	.056	.330	.327	.112	.940
	$\beta_3$	-.2	-.201	-.001	.207	.209	.043	.960	-.191	.009	.193	.203	.037	.969
	$\gamma_1$	-.1	-.096	.004	.169	.174	.029	.964	-.098	.002	.170	.177	.029	.965
	$\gamma_2$	.1	.098	-.002	.301	.302	.091	.948	.100	.000	.301	.314	.091	.954
	$\psi$	-.1	-.111	-.011	.316	.316	.100	.980	-.131	-.031	.404	.472	.164	.989
	$\sigma_b^2$	.5	.516	.016	.204	.217	.042	.949	.360	-.140	.138	.172	.038	.997
	$\Lambda(.9)$	.9	.917	.017	.191	.184	.061	.945	.932	.032	.196	.200	.052	.952
	$\Lambda(1.4)$	1.4	1.446	.046	.302	.297	.043	.945	1.471	.071	.314	.328	.040	.959
	$\Lambda(1.9)$	1.9	1.977	.077	.440	.442	.134	.958	2.013	.113	.460	.494	.122	.966
8	$\beta_0$	-1.0	-.994	.006	.201	.198	.040	.946	-.927	.073	.187	.188	.040	.928
	$\beta_1$	1.0	.990	-.010	.165	.168	.027	.954	.927	-.073	.154	.158	.029	.933
	$\beta_2$	-.5	-.504	-.004	.296	.288	.087	.948	-.471	.029	.277	.273	.078	.945
	$\beta_3$	-.2	-.206	-.006	.156	.156	.024	.953	-.199	.001	.148	.152	.022	.956
	$\gamma_1$	-.1	-.102	-.002	.179	.173	.032	.939	-.103	-.003	.179	.173	.032	.939
	$\gamma_2$	.1	.114	.014	.288	.300	.083	.962	.116	.016	.288	.305	.083	.965
	$\psi$	-.1	-.112	-.012	.230	.232	.053	.976	-.114	-.014	.266	.264	.071	.974
	$\sigma_b^2$	.5	.502	.002	.138	.142	.019	.961	.402	-.098	.106	.115	.021	.998
	$\Lambda(.9)$	.9	.906	.006	.176	.181	.040	.959	.913	.013	.178	.188	.035	.964
	$\Lambda(1.4)$	1.4	1.421	.021	.282	.289	.028	.958	1.432	.032	.287	.299	.025	.961
	$\Lambda(1.9)$	1.9	1.949	.049	.413	.428	.090	.965	1.964	.064	.419	.441	.081	.968
15	$\beta_0$	-1.0	-.989	.011	.168	.166	.028	.946	-.938	.062	.159	.159	.029	.931
	$\beta_1$	1.0	.997	-.003	.145	.142	.021	.943	.948	-.052	.137	.136	.021	.933
	$\beta_2$	-.5	-.508	-.008	.248	.245	.062	.950	-.481	.019	.235	.235	.055	.952
	$\beta_3$	-.2	-.201	-.001	.109	.113	.012	.958	-.198	.002	.105	.112	.011	.961
	$\gamma_1$	-.1	-.083	.017	.167	.172	.028	.955	-.084	.016	.168	.172	.028	.954
	$\gamma_2$	.1	.114	.014	.301	.299	.091	.951	.118	.018	.301	.301	.091	.953
	$\psi$	-.1	-.098	.002	.185	.190	.034	.956	-.090	.010	.200	.200	.040	.950
	$\sigma_b^2$	.5	.487	-.013	.098	.100	.010	.966	.424	-.076	.082	.085	.012	.944
	$\Lambda(.9)$	.9	.900	.000	.182	.179	.028	.955	.902	.002	.183	.182	.025	.957
	$\Lambda(1.4)$	1.4	1.400	.000	.297	.283	.021	.948	1.403	.003	.297	.287	.019	.953
	$\Lambda(1.9)$	1.9	1.914	.014	.431	.416	.062	.949	1.919	.019	.434	.422	.055	.952
25	$\beta_0$	-1.0	-.992	.008	.149	.150	.022	.949	-.951	.049	.142	.145	.023	.938
	$\beta_1$	1.0	.997	-.003	.132	.130	.017	.947	.957	-.043	.125	.125	.017	.934
	$\beta_2$	-.5	-.501	-.001	.223	.225	.050	.950	-.481	.019	.209	.217	.044	.964
	$\beta_3$	-.2	-.203	-.003	.086	.090	.007	.954	-.200	.000	.084	.089	.007	.960
	$\gamma_1$	-.1	-.101	-.001	.177	.172	.031	.941	-.098	.002	.174	.172	.030	.941
	$\gamma_2$	.1	.100	.000	.310	.299	.096	.943	.112	.012	.305	.300	.093	.947
	$\psi$	-.1	-.091	.009	.177	.169	.031	.947	-.084	.016	.185	.173	.034	.938
	$\sigma_b^2$	.5	.490	-.010	.083	.084	.007	.956	.446	-.054	.073	.073	.008	.931
	$\Lambda(.9)$	.9	.913	.013	.188	.182	.022	.941	.910	.010	.186	.183	.020	.944
	$\Lambda(1.4)$	1.4	1.428	.028	.305	.288	.018	.937	1.421	.021	.305	.289	.016	.932
	$\Lambda(1.9)$	1.9	1.958	.058	.454	.426	.053	.947	1.946	.046	.450	.426	.046	.946

Table 5.2: Summary of simulation results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) in the simultaneous modeling of binary longitudinal outcomes and survival time (n=400).

$n_i$	Par.	True	MLE						MPLE					
			Est.	Bias	SSD	ESE	MSE	CP	Est.	Bias	SSD	ESE	MSE	CP
4	$\beta_0$	-1.0	-1.003	-.003	.172	.170	.030	.949	-.924	.076	.158	.159	.031	.932
	$\beta_1$	1.0	1.004	.004	.145	.143	.021	.946	.929	-.071	.132	.133	.023	.916
	$\beta_2$	-.5	-.503	-.003	.248	.244	.061	.941	-.465	.035	.229	.229	.053	.944
	$\beta_3$	-.2	-.196	.004	.147	.147	.022	.952	-.186	.014	.137	.142	.019	.959
	$\gamma_1$	-.1	-.100	.000	.124	.121	.015	.936	-.101	-.001	.124	.121	.015	.938
	$\gamma_2$	.1	.114	.014	.213	.210	.046	.950	.116	.016	.213	.213	.046	.957
	$\psi$	-.1	-.096	.004	.205	.203	.042	.963	-.112	-.012	.262	.270	.069	.978
	$\sigma_b^2$	.5	.496	-.004	.136	.150	.018	.961	.349	-.151	.095	.119	.032	1.000
	$\Lambda(.9)$	.9	.898	-.002	.129	.126	.030	.948	.905	.005	.131	.131	.025	.950
	$\Lambda(1.4)$	1.4	1.406	.006	.208	.201	.021	.949	1.418	.018	.213	.208	.018	.953
	$\Lambda(1.9)$	1.9	1.916	.016	.302	.295	.062	.944	1.934	.034	.311	.306	.053	.951
8	$\beta_0$	-1.0	-.998	.002	.147	.139	.022	.938	-.930	.070	.138	.132	.024	.903
	$\beta_1$	1.0	1.001	.001	.120	.118	.014	.946	.937	-.063	.112	.111	.017	.904
	$\beta_2$	-.5	-.509	-.009	.205	.203	.042	.951	-.476	.024	.192	.192	.038	.944
	$\beta_3$	-.2	-.195	.005	.112	.109	.013	.939	-.188	.012	.107	.107	.012	.942
	$\gamma_1$	-.1	-.098	.002	.125	.120	.016	.953	-.098	.002	.125	.120	.016	.953
	$\gamma_2$	.1	.106	.006	.207	.209	.043	.948	.107	.007	.206	.210	.043	.950
	$\psi$	-.1	-.104	-.004	.155	.156	.024	.967	-.103	-.003	.178	.176	.032	.965
	$\sigma_b^2$	.5	.498	-.002	.093	.099	.009	.963	.401	-.099	.072	.080	.015	.937
	$\Lambda(.9)$	.9	.902	.002	.127	.126	.022	.943	.905	.005	.127	.128	.019	.942
	$\Lambda(1.4)$	1.4	1.413	.013	.206	.200	.015	.946	1.417	.017	.206	.203	.013	.949
	$\Lambda(1.9)$	1.9	1.924	.024	.299	.292	.043	.946	1.930	.030	.300	.296	.038	.949
15	$\beta_0$	-1.0	-.996	.004	.117	.117	.014	.946	-.944	.056	.110	.112	.015	.923
	$\beta_1$	1.0	1.000	.000	.099	.101	.010	.954	.949	-.051	.095	.096	.012	.925
	$\beta_2$	-.5	-.504	-.004	.172	.173	.030	.955	-.477	.023	.163	.166	.027	.955
	$\beta_3$	-.2	-.199	.001	.081	.080	.006	.943	-.196	.004	.078	.079	.006	.947
	$\gamma_1$	-.1	-.100	.000	.118	.120	.014	.962	-.100	.000	.118	.120	.014	.964
	$\gamma_2$	.1	.108	.008	.208	.209	.043	.955	.110	.010	.208	.209	.043	.954
	$\psi$	-.1	-.099	.001	.128	.130	.016	.959	-.089	.011	.140	.136	.020	.953
	$\sigma_b^2$	.5	.495	-.005	.070	.071	.005	.959	.431	-.069	.058	.060	.008	.886
	$\Lambda(.9)$	.9	.899	-.001	.126	.126	.014	.955	.901	.001	.126	.127	.012	.955
	$\Lambda(1.4)$	1.4	1.405	.005	.198	.199	.010	.951	1.408	.008	.198	.200	.009	.953
	$\Lambda(1.9)$	1.9	1.917	.017	.289	.290	.030	.948	1.915	.015	.286	.292	.027	.951
25	$\beta_0$	-1.0	-1.004	-.004	.104	.106	.011	.949	-.961	.039	.099	.102	.011	.939
	$\beta_1$	1.0	.998	-.002	.093	.092	.009	.945	.954	-.046	.087	.088	.010	.921
	$\beta_2$	-.5	-.486	.014	.158	.159	.025	.943	-.467	.033	.150	.153	.024	.944
	$\beta_3$	-.2	-.198	.002	.063	.063	.004	.962	-.197	.003	.063	.063	.004	.964
	$\gamma_1$	-.1	-.101	-.001	.118	.120	.014	.961	-.100	.000	.118	.120	.014	.965
	$\gamma_2$	.1	.099	-.001	.217	.208	.047	.941	.103	.003	.213	.209	.046	.948
	$\psi$	-.1	-.096	.004	.115	.117	.013	.959	-.081	.019	.120	.120	.015	.957
	$\sigma_b^2$	.5	.493	-.007	.058	.059	.003	.964	.446	-.054	.051	.052	.006	.872
	$\Lambda(.9)$	.9	.911	.011	.135	.127	.011	.931	.910	.010	.132	.127	.010	.940
	$\Lambda(1.4)$	1.4	1.419	.019	.213	.200	.009	.933	1.414	.014	.207	.200	.008	.940
	$\Lambda(1.9)$	1.9	1.925	.025	.299	.291	.026	.947	1.925	.025	.294	.292	.023	.949

gitudinal outcomes per subject; “True” gives the true values of parameters; the middle 6 columns under “MLE” and the right 6 columns under “MPLE” are the results of the maximum likelihood estimates from the EM algorithm and the proposed maximum penalized likelihood estimates, respectively; the averages of the estimates are in “Est.”; the averages of the bias estimates of the parameter estimates subtracted from true values are in “Bias”; the sample standard deviations from 1000 simulations are reported in “SSD”; “ESE” is the average of 1000 standard error estimates based on the observed information matrix; “MSE” gives the mean squared error calculated by adding the squared bias and the squared sample standard deviations; “CP” is the coverage proportion of 95% nominal confidence intervals based on the estimated standard error “ESE”. Note that “ESE” under “MPLE” is based on the observed information matrix obtained by maximum likelihood approach using the maximum penalized likelihood estimates. Satterthwaite method is used for the coverage proportion of  $\sigma_b^2$ .

From Table 5.1 and Table 5.2, we can see that the bias of the proposed MPLE is small for most cases although it is bigger than the MLE’s, but overall the bias of the MPLE decreases for the larger number of longitudinal observations per subject and the larger sample size like the MLE’s does. On the other hand, the estimate of  $\sigma_b^2$  of the MPLE is smaller than its true value showing the biggest bias, but it is improved soon being close to the true values as the number of longitudinal observations per subject increases. It is already known that the penalized quasi-likelihood (PQL) used for GLMMs tends to underestimate somewhat the variance components when applied to clustered binary data but the situation improves rapidly for binomial observations having denominators greater than one (Breslow and Clayton, 1993). The result from our simulation studies conforms this fact. For both MLE and MPLE, the estimated standard errors calculated from the observed information matrix are close to the sample standard deviations from the 1000 estimates. They decrease over the number of



longitudinal observations per subject except for the baseline cumulative hazards estimates, and they also decrease as sample size increases. The MPLE has smaller sample standard deviations and estimated standard errors than MLE for most cases. As for the mean squared error representing both bias and sample standard deviation together, the mean squared error of the MPLE appears to be smaller than or close to the MLE's. The mean squared errors from both MLE and MPLE decrease as the number of longitudinal observations per subject and sample size increase. The 95% confidence interval coverage rates are close to 0.95 except those for  $\psi$  of both MLE and MPLE with the small numbers of longitudinal observations per subject ( $n_i=4$  and 8) of the small sample size ( $n=200$ ), for  $\sigma_b^2$  of the MPLE with the small number of longitudinal observations per subject ( $n_i=4$  and 8) of the small sample size ( $n=200$ ), and for  $\sigma_b^2$  of the MPLE with the very small or large number of longitudinal observations per subject ( $n_i=4$ , 15 and 25) of the large sample size ( $n=400$ ). For both MLE and MPLE, the coverage rate of the parameter  $\psi$  is recovered for the large number of longitudinal observations per subject and the large sample size. Thus, with small number of longitudinal observations per subject and small sample size, the test for  $\psi$  is conservative, which strengthens the test results when rejecting the null ( $\psi = 0$ ), and the type I error becomes closer to the nominal level as the number of longitudinal observations per subject and sample size increase. While the high coverage rate of  $\sigma_b^2$  of the MLE is improved for both large number of longitudinal observations per subject and large sample size, the coverage rate of  $\sigma_b^2$  of the MPLE appears to be improved for the large number of longitudinal observations per subject with small sample size and the small number of longitudinal observations per subject with the large sample size. With the small sample size of 200 of Table 5.1, the high coverage rates of  $\sigma_b^2$  of the MPLE at the small numbers of longitudinal observations per subject ( $n_i=4$  and 8) are recovered at the large numbers of longitudinal observations per subject ( $n_i=15$  and 25). On the other hand, with the

relatively small number of longitudinal observations per subject ( $n_i=8$ ), the high coverage rate of  $\sigma_b^2$  of the MPLE shown at the small sample size of 200 in Table 5.1 is improved for the large sample size of 400 in Table 5.1. In additional simulation studies conducted with the larger sample size of 800 whose results are not provided in this paper, the high coverage rates of  $\sigma_b^2$  of the MPLE shown at the smallest number of longitudinal observations per subject ( $n_i=4$ ) with sample sizes of 200 and 400 in both Tables 5.1 and 5.2 actually reached 95% nominal level for the sample size of 800. Figure 5.1 shows the ratios of mean squared errors (MSEs) of the proposed MPLE to the MLE with sample sizes of 200 and 400 for the parameters of predictors in longitudinal and hazard models. This figure confirms the results provided in Table 5.1 and Table 5.2 in that all plots indicate the ratios of mean squared errors are close to 1 which implies the proposed MPLE provides the mean squared errors close to the MLE's. Figure 5.2 shows the ratios of user times of the proposed MPLE to the MLE with sample sizes of 200 and 400. Both plots show the proposed MPLE is more efficient reducing about 70% of the computing time of the MLE over all different numbers of longitudinal outcomes per subject and sample sizes.

## 5.5 Analysis of the CHANCE Study

The Carolina Head and Neck Cancer Study (CHANCE) is a population based epidemiologic study conducted at 60 hospitals in 46 counties in North Carolina from 2002 through 2006 (Divaris *et al.* 2010). Patients were diagnosed with head and neck cancer (oral, pharynx, and larynx cancer) from 2002–2006. Their survival status was collected up to 2007 and QoL was evaluated over time for three years after diagnosis. QoL information was collected through questionnaires. Based on summary scores of the five domains of self-perceived quality of life including Physical Well-Being (PWB), Social/Family Well-Being (SWB), Emotional Well-Being (EWB), Functional Well-Being

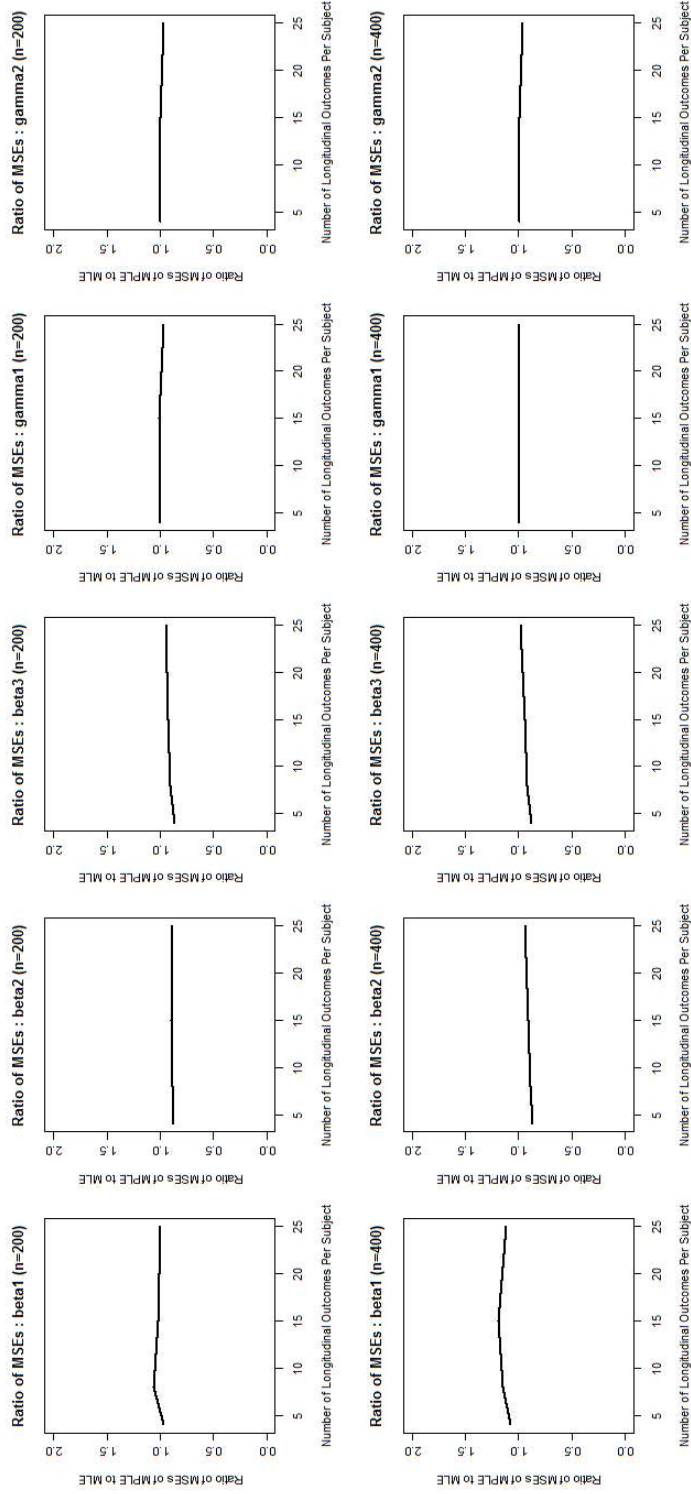


Figure 5.1: Plot of ratios of mean squared errors (MSEs) of maximum penalized likelihood estimator (MPLE) to maximum likelihood estimator (MLE) for parameters of predictors in longitudinal and hazard models (n=200, 400)

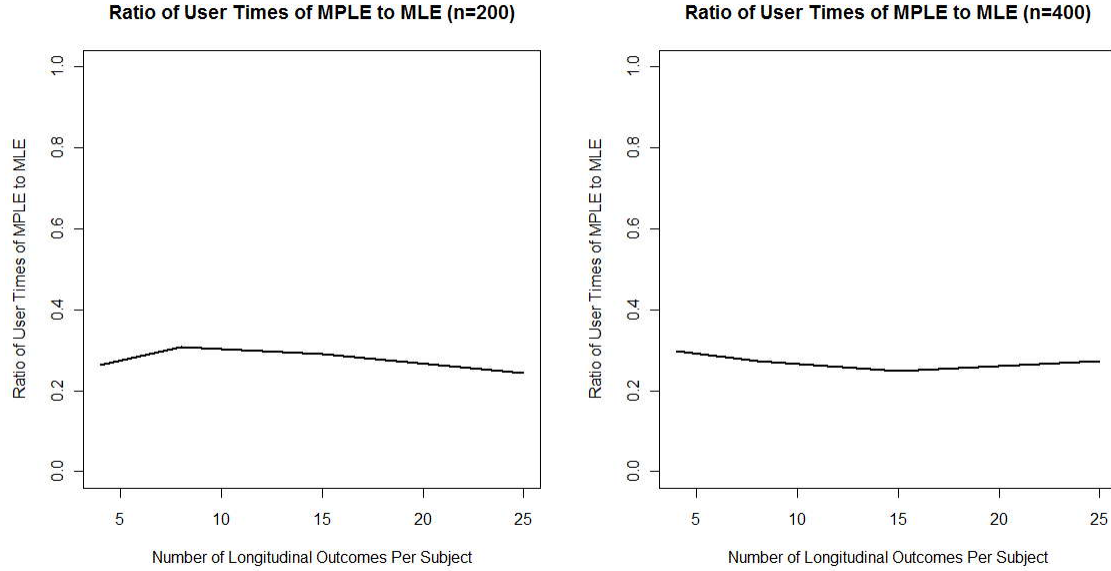


Figure 5.2: Plot of ratios of user times of maximum penalized likelihood estimator (MPLE) to maximum likelihood estimator (MLE) ( $n=200, 400$ )

(FWB) and Head and Neck Cancer Specific symptoms (HNCS), patient's QoL information was classified into satisfaction or dissatisfaction with life. Survival time is defined as the time to death from diagnosis. Demographic and life style characteristics, medical histories and clinical factors are also collected. Ending in December 2009 and excluding the patients with missing data, information on QoL has been obtained from 554 head and neck cancer patients in the analysis. Based on the death information through 2007 available from the National Death Index (NDI), 85 of 554 patients died and the censoring rate is 85%. The number of observations per patient ranges 1 to 3 with average of 1.93. It is of interest to elucidate the variables which are associated with both QoL satisfaction and survival time for patients with head and neck cancer. In particular, we are interested in the comparison between African-Americans and Whites since it is known that African-Americans have a higher incidence of head and neck cancer and worse survival than Whites. The longitudinal QoL satisfaction outcomes and survival time are correlated within a patient, and this dependency should be taken into account

in the analysis.

We apply both approaches of the maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) to Head and Neck Cancer Specific symptoms (HNCS) among QoL domains with survival time. Longitudinal HNCS QoL outcomes are binary measurements with 1 (“satisfied”) and 0 (“dissatisfied”). In both longitudinal QoL and hazards models, we consider race (African-Americans, Whites), the number of 12 oz. beers consumed per week (None, <1, 1–4, 5–14, 15–29,  $\geq 30$ ), household income (0–10K, 20–30K, 40–50K,  $\leq 60$ K), surgery (Yes/No), radiation therapy (Yes/No), chemotherapy (Yes/No), primary tumor site (Oral & Pharyngeal, Laryngeal) and tumor stage (I, II, III, IV) as categorical, and age at diagnosis (range: 24–80), the number of persons supported by household income (range: 1–5), body mass index (BMI) (range: 15.66–56.28) and the total number of medical conditions reported (range: 0–6) as continuous. Additionally, 2 interactions with race, i.e. race  $\times$  the total number of medical conditions reported and race  $\times$  tumor site, are included in both models since we are particularly interested in the difference of QoL and survival between African American and White. Time at survey measurement is also included as a covariate for longitudinal outcomes. A random intercept for the dependence between the QoL satisfaction and the risk of death is included in both models, and assumed to follow a mixture of normal distributions. In Table 5.3, we compare the estimates and the estimated standard errors of the maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE). From both “Est.” and “ESE” columns, we see that MLE and MPLE provide similar estimates and estimated standard errors each other for the parameters of interest in longitudinal QoL and hazards models. On the other hand, the parameters of  $\sigma_b^2$  and  $\psi$ , which denote the variance of random effects and the coefficient of random effects characterizing the dependence between longitudinal QoL and survival processes, respectively, have different estimates and estimated

Table 5.3: Analyses results from maximum likelihood estimation (MLE) and maximum penalized likelihood estimation (MPLE) for the Quality of Life and survival time for the CHANCE study

Parameter		Est.		ESE	
		MLE	MPLE	MLE	MPLE
<i>HNCS QoL longitudinal model</i>					
Intercept	$\beta_0$	.197	.206	.824	.922
Race (ref= White): African American	$\beta_1$	.562	.564	.345	.391
# of 12 oz. beers consumed per week (ref=30 or more)					
– None	$\beta_2$	.615	.618	.275	.315
– less than 1	$\beta_3$	.751	.706	.366	.409
– 1 to 4	$\beta_4$	1.259	1.253	.300	.337
– 5 to 14	$\beta_5$	1.062	1.072	.257	.294
– 15 to 29	$\beta_6$	.581	.577	.294	.336
Household income (ref= level1: 0–10K)					
– level2: 20–30K	$\beta_7$	- .268	- .254	.236	.271
– level3: 40–50K	$\beta_8$	.291	.309	.257	.293
– level4: $\geq$ 60K	$\beta_9$	1.181	1.162	.279	.313
Surgery (ref= No): Yes	$\beta_{10}$	- .029	- .032	.205	.234
Radiation therapy (ref= No): Yes	$\beta_{11}$	-1.179	-1.098	.299	.323
Chemotherapy (ref= No): Yes	$\beta_{12}$	.219	.197	.245	.280
Tumor site (ref=Oral & Pharyngeal): Laryngeal	$\beta_{13}$	- .225	- .218	.225	.255
Tumor stage (ref= I)					
– II	$\beta_{14}$	- .479	- .470	.306	.334
– III	$\beta_{15}$	-1.494	-1.418	.320	.358
– IV	$\beta_{16}$	-1.383	-1.342	.308	.343
Age at diagnosis	$\beta_{17}$	.012	.012	.009	.011
# of persons supported by household income	$\beta_{18}$	- .174	- .176	.088	.100
BMI	$\beta_{19}$	.042	.040	.015	.017
Total # of medical conditions reported	$\beta_{20}$	- .208	- .203	.092	.104
Race (African-American) $\times$ Tumor site (Laryngeal)	$\beta_{21}$	- .156	- .178	.438	.496
Race (African-American) $\times$ Total # of medical conditions reported	$\beta_{22}$	.088	.095	.197	.224
Time at survey measurement (years)	$\beta_{23}$	.243	.216	.067	.070
variance of random effects	$\sigma_b^2$	.317	1.037	.185	.394
<i>Hazards model</i>					
Random effect coefficient	$\psi$	-1.560	- .623	1.060	.285
Race (ref= White): African American	$\gamma_1$	.482	.411	.461	.384
# of 12 oz. beers consumed per week (ref=30 or more)					
– None	$\gamma_2$	- .866	- .795	.417	.354
– less than 1	$\gamma_3$	- .241	- .198	.444	.395
– 1 to 4	$\gamma_4$	- .915	- .845	.447	.389
– 5 to 14	$\gamma_5$	-1.180	-1.106	.409	.350
– 15 to 29	$\gamma_6$	- .616	- .568	.422	.372
Household income (ref= level1: 0–10K)					
– level2: 20–30K	$\gamma_7$	- .168	- .195	.321	.283
– level3: 40–50K	$\gamma_8$	- .898	- .852	.400	.353
– level4: $\geq$ 60K	$\gamma_9$	-1.359	-1.319	.446	.396
Surgery (ref= No): Yes	$\gamma_{10}$	- .489	- .501	.324	.282
Radiation therapy (ref= No): Yes	$\gamma_{11}$	- .454	- .468	.411	.361
Chemotherapy (ref= No): Yes	$\gamma_{12}$	.065	.052	.372	.334
Tumor site (ref=Oral & Pharyngeal): Laryngeal	$\gamma_{13}$	- .012	- .018	.333	.295
Tumor stage (ref= I)					
– II	$\gamma_{14}$	- .163	- .240	.489	.428
– III	$\gamma_{15}$	.302	.179	.506	.426
– IV	$\gamma_{16}$	1.239	1.086	.478	.375
Age at diagnosis	$\gamma_{17}$	.023	.017	.019	.008
# of persons supported by household income	$\gamma_{18}$	.086	.059	.137	.110
BMI	$\gamma_{19}$	.015	.011	.022	.016
Total # of medical conditions reported	$\gamma_{20}$	.277	.249	.134	.112
Race (African-American) $\times$ Tumor site (Laryngeal)	$\gamma_{21}$	.342	.339	.573	.505
Race (African-American) $\times$ Total # of medical conditions reported	$\gamma_{22}$	- .256	- .242	.270	.245

standard errors between the MLE and MPLE. This discrepancy of the MPLE from the MLE may be a numerical issue due to the small cluster size with the average of 1.93. In addition, the MPLE provides slightly bigger estimated standard errors than the MLE for the parameters in the longitudinal model while it appears in the reverse direction in hazards model. This also may be a numerical issue due to the small number of longitudinal outcomes per subject since the estimation in the longitudinal model is directly affected by the individual cluster size while the estimation in hazards model is not. Comparing the computing time spent on producing the results in Table 5.3, the proposed MPLE took only a sixth of the time the MLE did (62.83 and 361.78 seconds for MPLE and MLE respectively). This analysis result indicates that, even for the small cluster size, the proposed MPLE provides the similar results to those of the MLE for the parameters of interest taking less computing time than the MLE. In the studies with larger number of longitudinal outcomes per subject, the results of the MPLE are expected to be close to those of the MLE for all parameters with much better efficiency on calculation.

## 5.6 Concluding Remarks

In this paper, we have developed a more computationally efficient estimation procedure adopting a penalized likelihood based on Laplace approximation for the simultaneous modeling of longitudinal outcomes and survival time. Our proposed penalized likelihood estimation method is an effort to reduce the intensity on computation still providing the similar estimates to those by the EM algorithm of the maximum likelihood approach. Simulation studies indicated that the penalized likelihood approach performs as well as the EM algorithm of maximum likelihood approach, but only requires a fraction of the computing time. We also illustrated this comparison with the CHANCE data.

For the purpose of comparison, we also conducted the simulation studies of two ad-

ditional approximation methods which are using the Laplace approximation to the full log-likelihood in (5.10) and using the observed variance of the longitudinal outcomes and the observed event for the expected values in the penalized log-likelihood (5.18) for binary longitudinal outcomes and survival time. The simulation results indicated that the former provides the estimates more close to the maximum likelihood estimates but takes longer time than the two penalized likelihood procedures although its computing time is also less than the maximum likelihood approach. It is because, from the Laplace approximated full log-likelihood, all parameters in  $|\mathbf{I}_{db} - \Sigma_b \tilde{l}''_{i|b}(\boldsymbol{\theta}, \Lambda_s)|$  are used for estimation and thus the calculation is more complicated than the penalized likelihood estimation methods. In the mean time, the latter provides the most biased estimates but takes the least computing time since it uses the most approximated expression of log-likelihood. Therefore, in terms of bias and computing time, the proposed penalized likelihood method using the original expressions for the expected values in the penalized likelihood appears to behave best among the three approximation methods.

In the simultaneous modeling considered in this paper, we assumed random effects to follow a Gaussian distribution with mean zero. However, it is unclear whether the normality assumption is truly satisfied in practice. Future work can include developing an approach to diminish computational intensity efficiently through the penalized likelihood approach for relaxing the normality assumption of random effects in the simultaneous modeling.



# Chapter 6

## SUMMARY AND FUTURE RESEARCH

In this dissertation, we have studied statistical methods for joint analysis of survival time and longitudinal data and particularly proposed the simultaneous modeling of the two different types of outcomes. Random effects are introduced to account for the dependence between longitudinal outcomes and survival time due to unobserved factors. Specifically, in terms of the distributional assumption of random effects, the following two scenarios were studied: 1) assuming a Gaussian process for random effects, 2) assuming the underlying distribution of random effects to be unknown.

This dissertation research is motivated by biomedical and public health applications where it is common that both longitudinal outcomes over time and survival endpoint are collected for the same subject along with the subject's characteristics or risk factors. Investigators are often interested in finding important variables for predicting both longitudinal outcomes and survival time which are correlated within a subject. This naturally led us to consider simultaneous models. Thus, the main contribution of this dissertation is to provide statistical methods which address the association of covariates to both longitudinal outcomes and survival time of interest while the dependence between the two outcomes is taken into account. Generalized linear mixed model was considered for longitudinal process in order to incorporate both categorical

and continuous longitudinal outcomes. A stratified proportional hazards model was assumed for survival time. The cumulative baseline hazard functions were also studied and Breslow-type estimates were proposed.

In Chapter 3, we have assumed random effects to follow a multivariate Gaussian process in the joint analysis and considered a nonparametric maximum likelihood estimation approach. The EM algorithm was adopted for estimation, and the proposed estimates performed well in finite samples under the various simulation settings considered. The variance estimates based on the observed information matrix approximated the true variance well in finite samples. In Chapter 4, we have relaxed normality assumption for random effects in the joint analysis. Assuming the underlying distribution of random effects to be unknown, we used a mixture of Gaussian distributions as an approximation for the random effect distribution. We also developed a nonparametric maximum likelihood estimation method, and weights of the mixture components were estimated with model parameters using the EM algorithm. The proposed estimators were shown to have nice finite sample properties via simulation studies. The simulation studies conducted for robustness of the assumed mixture distribution indicated that, when the true distribution of random effects is not normal, the mixture of normal distributions well-approximate the random effect distribution by yielding the less biased estimates for the parameters of interest in longitudinal and hazards models and the more similar shaped density plot to the true distribution than no mixture. The number of mixture distributions was shown to be properly selected by AIC and BIC through simulation studies. For both maximum likelihood approaches with and without normality assumption of random effects studied in Chapters 3 and 4 respectively, the proposed estimators were proved to have desirable asymptotic properties such as consistency and asymptotic normality, and most of the proofs relied on the empirical processes theory. In Chapter 5, we have considered the penalized likelihood for

a more computationally efficient estimation procedure than the EM algorithm of the maximum likelihood approach with a Gaussian process for random effects. Simulation studies showed the proposed penalized likelihood approach reduced the computational intensity spending only a fraction of the computing time the EM algorithm took, still providing the similar estimates and mean squared errors to those by the EM algorithm.

All proposed methods in this dissertation were illustrated with the real-world data sets from the CHANCE study. We compared the results from the proposed simultaneous analysis and separate analyses in Chapters 3 and 4. Under both situations of normally distributed and distributional free random effects, the simultaneous analysis additionally identified more predictors for longitudinal quality of life and survival time than separate analyses. In Chapter 5, we compared the analysis results from the maximum likelihood approach and maximum penalized likelihood approach for assuming normality of random effects. The latter showed the remarkable reduction from the former's computing time while both produced the similar results in estimation.

The proposed methods in this dissertation research can be extended in several directions:

First, in this dissertation, we considered the generalized linear mixed model for longitudinal process to incorporate both categorical and continuous longitudinal outcomes. In the the joint analysis framework, relatively little work was done for categorical longitudinal data while continuous longitudinal data were studied by many authors. Our proposed approaches generalize previous work to general longitudinal outcomes and this work fills in some gaps in the joint modeling research. Then, future work can include considering generalization to mixed types of longitudinal outcomes.

Second, in some applications where sample size and the number of observations per subject are too large, the EM algorithm of maximum likelihood approach may be intensive on computation due to the integration of complete data likelihood over

random effects. Thus, in my dissertation, we considered a penalized likelihood for the simultaneous modeling with a Gaussian process of random effects. It will be also of interest to develop an approach to relieve computational burden efficiently through the penalized likelihood approach for relaxing the normality assumption of random effects in the simultaneous modeling.

Third, we considered one survival event in the simultaneous analysis proposed in this dissertation, but one may be interested in multivariate endpoints such as recurrent events, multiple disease outcome data and competing risk factors. Therefore, exploring the possibility of the extension of the proposed approaches to multivariate survival data would be worth pursuing.

Last, but not least, in real applications, the proportional hazards assumption considered in this dissertation may not always be true or one may be interested in modeling association from different aspects. A natural extension would be to consider other types of models for survival time including, but not limited to the proportional odds model, the accelerated failure time model, or the additive hazards model.

# REFERENCES

- Agresti, A., Caffo, B., and Ohman-Strickland, P. (2004), Examples in Which Misspecification of a Random Effects Distribution Reduces Efficiency, and Possible Remedies, *Computational Statistics and Data Analysis*, **47**, 639–653.
- Albert, P.S. and Follmann, D.A. (2000), Modeling Repeated Count Data Subject to Informative Dropout, *Biometrics*, **56**, 667–677.
- Albert, P.S. and Follmann, D.A. (2007), Random Effects and Latent Processes Approaches for Analyzing Binary Longitudinal Data with Missingness: a Comparison of Approaches Using Opiate Clinical Trial Data, *Statistical Methods in Medical Research*, **16**, 417–439.
- Albert, P.S., Follmann, D.A., Wang, S.A., and Suh, E.B. (2002), A Latent Autoregressive Model for Longitudinal Binary Data subject to Informative Missingness, *Biometrics*, **58**, 631–642.
- Andersen, P.K. and Gill, R.D. (1982), Cox’s Regression Model for Counting Processes: A Large Sample Study, *The Annals of Statistics*, **10**, 1100–1120.
- Andersen, P.K, Borgan, O., Gill, R.D. and Keiding, N. (1993), *Statistical Models Based on Counting Processes*, Springer–Verlag.
- Barndorff-Nielsen, O.E. and Cox, D.R. (1989), *Asymptotic Techniques for Use in Statistics*, London: Chapman and Hall.
- Bartholomew, D.J. (1987), *Latent Variable Models and Factor Analysis*, Oxford University Press, New York.
- Bartlett, R.F. and Sutradhar, B.C. (1999), On estimating equations for parameters in generalized linear mixed models with application to binary data, *Environmetrics*, **10**, 769–784.
- Bickel, P.J., Klaassen, C.A.J., Ritov, Y., and Wellner, J. A. (1993), *Efficient and Adaptive Estimation for Semiparametric Models*, Johns Hopkins University Press,

Baltimore.

- Böhning, D. (1999), *Computer-Assisted Analysis of Mixtures and Applications: Meta-analysis, Disease Mapping and Others*, Number 81 in Monographs on Statistics and Applied Probability, Chapman & Hall/CRC.
- Breslow, N.E. and Clayton, D.G., (1993), Approximate Inference in Generalized Linear Mixed Models, *Journal of American Statistical Association*, **88**, 125-134.
- Breslow, N.E. and Day, N.E. (1980), *Statistical Methods in Cancer Research, Volume I*. IARC Scientific Publications No. 32. Lyon.
- Breslow, N.E. and Lin, X. (1995), Bias Correction in Generalized Linear Mixed Models with a Single Component of Dispersion, *Biometrika*, **82**, 81-91.
- Brown, E.R. and Ibrahim, J.G. (2003), A Bayesian Semiparametric Joint Hierarchical Model for Longitudinal and Survival Data, *Biometrics*, **59**, 221-228.
- Buzas, J.S. (1998), Unbiased Scores in Proportional Hazards Regression with Covariate Measurement Error, *Journal of Statistical Planning and Inference*, **67**, 247-257.
- Cai, J. and Prentice, R.L (1995), Estimating Equations for Hazard Ratio Parameters Based on Correlated Failure Time Data, *Biometrika*, **82**, 151-164.
- Cai, J. and Prentice, R.L. (1997), Regression Analysis for Correlated Failure Time Data, *Lifetime Data Analysis*, **3**, 197-213.
- Caffo, B., Ming-Wen, A., and Rohde, C. (2007), Flexible Random Intercept Models for Binary Outcomes Using Mixtures of Normals, *Computational Statistics and Data Analysis*, **51**, 5220-5235.
- Carroll, R.J., Ruppert, D. and Stefanski, L.A. (1995), *Measurement Error in Nonlinear Models*, London: Chapman and Hall.
- Cella, D.F. (1994), *Manual: Functional Assessment of Cancer Therapy (FACT) Scales and Functional Assessment of HIV Infection (FAHI) Scale*, Chicago: Rush-Presbyterian-St. Luke's Medical Center
- Cella, D.F., Tulsky, D.S., Gray, G., Sarafian, B., Linn, E., Bonomi, A., Silberman, M., Yellen, S.B., Winicour, P., and Brannon, J. (1993), The Functional Assessment

of Cancer Therapy Scale: Development and Validation of the General Measure, *Journal of Clinical Oncology*, **11**, 570–579

Clayton, D. and Cuzick, J. (1985), Multivariate Generalizations of the Proportional Hazards Model (with discussion), *Journal Royal Statistical Society A*, **148**, 82–117.

Clegg, L.X, Cai, J. and Sen, P.K. (1999), A Marginal Mixed Baseline Hazards Model for Multivariate Failure Time Data, *Biometrics*, **55**, 805–812.

Cox, D.R. (1972), Regression Models and Life-tables (with discussion), *Journal Royal Statistical Society B*, **34**, 187–220.

Cox, D.R. (1975), Partial Likelihood, *Biometrika*, **62**, 269–276.

Cox, D.R. and Oakes, D. (1984), *Analysis of Survival Data*, London: Chapman & Hall.

Crouch, A.C. and Spiegelman, E. (1990), The Evaluation of Integrals of the Form  $\int f(t) \exp(-t^2) dt$ : Application to Logistic-normal Models, *Journal of American Statistical Association*, **85**, 464–469.

Dang, Q.Y., Mazumdar, S., and Houck, P.R. (2008), Sample size and power calculations based on generalized linear mixed models with correlated binary outcomes, *Computer Methods and Programs in Biomedicine*, **91**, 122–127.

D’Antonio, L.L., Zimmerman, G.J., Cella, D.F., and Long, S.A. (1996), Quality of Life and Functional Status Measures in Patients with Head and Neck Cancer, *Archives of Otolaryngology Head & Neck Surgery*, **122**, 482–487

Davidian, M. and Gallant, A.R. (1992), The Nonlinear Mixed Effects Model with a Smooth Random Effects Density. *Department of Statistics Technical Report, North Carolina State University*, Campus Box 8203, Raleigh, North Carolina 27695.

DeGruttola, V. and Tu, X.M. (1994), Modeling Progression of CD-4 Lymphocyte Count and Its Relationship to Survival Time, *Biometrics*, **50**, 1003–1014.

Dempster, A.P., Laird, N.M., and Rubin, D.B. (1977), Maximum Likelihood from Incomplete Data via the EM algorithm, *Journal Royal Statistical Society B*, **39**, 1–38.

Diggle, P.J. (1988), An Approach to the Analysis of Repeated Measurements, *Bio-*

*metrics*, **44**, 959–971.

Diggle, P.J., Heagerty, P., Liang, K.Y. and Zeger, S.L. (2002), *Analysis of Longitudinal Data*, Oxford University Press 2nd Ed.

Ding, J. and Wang, J.L. (2008), Modeling Longitudinal Data with Nonparametric Multiplicative Random Effects Jointly with Survival Data, *Biometrics*, **64**, 546–556.

Divaris, K., Olshan, A.F., Smith, J., Bell, M.E., Weissler, M.C., Funkhouser, W.K., and Bradshaw, P.T. (2010), “Oral Health and Risk for Head and Neck Squamous Cell Carcinoma: the Carolina Head and Neck Cancer Study, *Cancer Causes Control*, **21**, 567–575.

Drum, M.L. and McCullagh, P. (1993), REML Estimation with Exact Covariance in the Logistic Mixed Model, *Biometrics*, **49**, 677–689.

Elashoff, R.M., Li, G., and Li, N. (2007), An Approach to Joint Analysis of Longitudinal Measurements and Competing Risks Failure Time Data, *Statistics in Medicine*, **26**, 2813–2835.

———(2008), A Joint Model for Longitudinal Measurements and Survival Data in the Presence of Multiple Failure Types, *Biometrics*, **64**, 762–771.

Fang, F.M., Chien, C.Y., Kuo, S.C., Chiu, H.C., and Wang, C.J. (2004), Changes in quality of life of head-and-neck cancer patients following postoperative radiotherapy, *Acta Oncologica*, **43**, 571–578.

Faucett, C.L., Schenker, N., and Elashoff, R.M. (1998), Analysis of Censored Survival Data with Intermittently Observed Time-Dependent Binary Covariates, *Journal of American Statistical Association*, **93**, 427–437.

Faucett, C.J. and Thomas, D.C. (1996), Simultaneously Modeling Censored Survival Data and Repeatedly Measured Covariates: A Gibbs Sampling Approach, *Statistics in Medicine*, **15**, 1663–1685.

Fieuws, S., Spiessens, B., and Draney, K. (2004), *Mixture Models*. In: De Boeck, P., Wilson, M. (Eds.), *Explanatory Item Response Models: A Generalized Linear and Nonlinear Approach*, Springer-Verlag, New York, Ch. 11, 317–340.



- Firth, D. (1991), Generalized Linear Models, *Statistical Theory and Modelling*, eds. Hinkley, D.V. and Snell, E.J, London: Chapman and Hall, 55–82.
- Flemming, T.R. and Harrington, D.P. (1991), *Counting Processes and Survival Analysis*, New York: Wiley.
- Gallant, A.R. and Nychka, D.W. (1987), Semi-nonparametric Maximum Likelihood Estimation, *Econometrica*, **55**, 363–390.
- Ghidey, W., Lesaffre, E., and Eilers, P. (2004), Smooth Random Effects Distribution in a Linear Mixed Model, *Biometrics*, **60**, 945–953.
- Gilmour, A.R., Anderson, R.D. and Rae, A.L. (1985), The Analysis of Binomial Data by a Generalized Linear Mixed Model, *Biometrika*, **72**, 593–599.
- Green, P.J. (1987), Penalized Likelihood for General Semi-parametric Regression Models, *International Statistical Review*, **55**, 245–259.
- Heagerty, P.J. and Kurland, B.F. (2001), Misspecified Maximum Likelihood Estimates and Generalised Linear Mixed Models, *Biometrika*, **88**, 973–985.
- Henderson, R., Diggle, P., and Dobson, A. (2000), Joint Modeling of Longitudinal Measurements and Event Time Data, *Biometrics*, **4**, 465–480.
- Hogan, J. and Laird, N. (1997), Mixture Models for the Joint Distribution of Repeated Measures and Event Times, *Statistics in Medicine*, **16**, 239–257.
- Holloway, R.L., Hellewell, J.L., Marbella, A.M., Layde, P.M., Myers, K.B., and Campbell, B.H. (2005), Psychosocial effects in long-term head and neck cancer survivors, *Head and Neck – Journal for the Sciences and Specialties of the Head and Neck*, **27**, 281–288.
- Hougaard, P. (2000), *Analysis of Multivariate Survival Data*, New York: Springer.
- Hsieh, F., Tseng, Y.K. and Wang, J.L. (2006), Joint Modeling of Survival and Longitudinal Data: Likelihood Approach Revisited, *Biometrics*, **62**, 1037–1043.
- Hu, W., Li, G., and Li, N. (2009), A Bayesian Approach to Joint Analysis of Longitudinal Measurements and Competing Risks Failure Time Data, *Statistics in Medicine*, **28**, 1601–1619.

- Huang, W., Zeger, S., Anthony J., and Garrett E. (2001), Latent Variable Model for Joint Analysis of Multiple Repeated Measures and Bivariate Event times, *Journal of American Statistical Association*, **96**, 906–914.
- Huber, P., Ronchetti, E., and Victoria-Feser, M.P. (2004), “Estimation of generalized linear latent variable models”, *Journal Royal Statistical Society B*, **66**, 893–908.
- Jang, W. and Lim, J. (2009), A Numerical Study of PQL Estimation Biases in Generalized Linear Mixed Models Under Heterogeneity of Random Effects, *Communications in Statistics-Simulation and Computation*, **38**, 692–702.
- Kalbfleish, J.D. and Prentice, R.L. (2002), *The Statistical Analysis of Failure Time Data*, Wiley, John & Sons 2nd Ed.
- Klein, J.P. (1992), Semiparametric Estimation of Random Effects Using the Cox Model Based on the EM Algorithm, *Biometrics*, 48, 795–806.
- Kleinman, K.P. and Ibrahim, J.G. (1998), A Semiparametric Bayesian Approach to the Random Effects Model, *Biometrics*, **54**, 921–938.
- Komárek, A. and Lesaffre E. (2008a), “Generalized Linear Mixed Model with a Penalized Gaussian Mixture as a Random Effects Distribution, *Computational Statistics and Data Analysis*, 52, 3441–3458.
- (2008b), Bayesian Accelerated Failure Time Model with Multivariate Doubly Interval-Censored Data and Flexible Distributional Assumptions, *Journal of American Statistical Association*, **103**, 523–533.
- (2009), The Regression Analysis of Correlated Interval-censored Data: Illustration Using Accelerated Failure Time Models with Flexible Distributional Assumptions, *Statistical Modelling*, **9**, 299–319.
- Laird N.M. (1978), Empirical Bayes Methods for Two-Way Contingency Tables, *Biometrika*, **65**, 581–590.
- Laird, N. (1978), Nonparametric Maximum Likelihood Estimation of a Mixing Distribution, *Journal of American Statistical Association*, **73**, 805–811.
- Laird, N.M. and Ware, J.H. (1982), Random Effects Models for Longitudinal Data, *Journal of American Statistical Association*, **73**, 963–974.

- Lange, N. and Ryan, L. (1989), Assessing Normality in Random Effects Models, *The Annals of Statistics*, **17**, 624–642
- Larsen K. (2004), Joint Analysis of Time-to-Event and Multiple Binary Indicators of Latent Classes, *Biometrics*, **60**, 85–92.
- Lee, E.W., Wei, L.J., and Amato, D.A. (1992), Cox-Type Regression Analysis for Large Numbers of Small Groups of Correlated Failure Time Observations, *Survival Analysis: State of the Art. J. P. Klein and P.K. Goel (eds.)*, Kluwer Academic Publishers, 237–247.
- Li, Y. and Lin, X. (2000), Covariate Measurement Errors in Frailty Models for Clustered Survival Data, *Biometrika*, **87**, 846–866.
- Li, E., Wang, N. and Wang, N.Y. (2007), Joint Models for a Primary Endpoint and Multiple Longitudinal Covariate Processes, *Biometrics*, **63**, 1068–1078.
- Liang, K.Y. and Zeger, S.L. (1986), Longitudinal Data Analysis Using Generalized Linear Models, *Biometrika*, **73**, 13–22.
- Liang, K.Y., Self, S.G., and Chang, Y.(1993), Modeling Marginal Hazards in Multivariate Failure Time Data, *Journal Royal Statistical Society B*, **55**, 441–453.
- Lin, D.Y. (1994), Cox Regression Analysis of Multivariate Failure Time Data: The Marginal Approach, *Statistics in Medicine*, **13**, 2233–2247.
- Lin, D.Y. (2000), On Fitting Cox’s Proportional Hazards Models to Survey Data, *Biometrika*, **87**, 37–47.
- Lin, X. and Breslow, N.E. (1996), Bias Correction in Generalized Linear Mixed Models with Multiple Components of Dispersion, *Journal of American Statistical Association*, **91**, 1007–1016.
- List, L.L., D’Antonio, L.L., Cella, D.F., Siston, A., Mumby, P., Haraf, D., and Vokes, E. (1996), The Performance Status Scale for Head and Neck Cancer Patients and the Functional Assessment of Cancer Therapy–Head and Neck (FACT–H&N) Scale: A Study of Utility and Validity, *Cancer*, **77**, 2294–2301
- Localio, A.R., Berlin, J.A., and Ten Have T.R. (2006), Longitudinal and repeated cross-sectional cluster-randomization designs using mixed effects regression for binary

- outcomes: Bias and coverage of frequentist and Bayesian methods, *Statistics in Medicine*, **25**, 2720–2736.
- Louis, T.A. (1982), Finding the Observed Information Matrix when Using the EM Algorithm, *Journal Royal Statistical Society B*, **44**, 226–233
- Masaoud, E. and Stryhn, H. (2010), A simulation study to assess statistical methods for binary repeated measures data, *Preventive Veterinary Medicine*, **93**, 81–97.
- McDullagh, P. and Nelder, J.A. (1989), *Generalized Linear Models*, Chapman & Hall/CRC 2nd Ed.
- Nakamura, T. (1992), Proportional Hazards Models with Covariates Subject to Measurement Error, *Biometrika*, **48**, 829–838.
- Nelson, K.P. and Leroux, B.G. (2008), Properties and comparison of estimation methods in a log-linear generalized linear mixed model, *Journal of Statistical Computation and Simulation*, **78**, 367–384.
- Neuhaus, J.M., Hauck, W.W., and Kalbfleisch, J.D. (1992), The Effects of Mixture Distribution Misspecification When Fitting Mixed-Effects Logistic Models, *Biometrika*, **79**, 755–762.
- Neyman, J. and Scott, E.L. (1948), Consistent estimates based on partially consistent observations, *Econometrica*, **16**, 1–32.
- Nibu, K.I., Ebihara, Y., Ebihara, M., Kawabata, K., Onitsuka, T., Fujii, T., and Saikawa, M. (2010), Quality of life after neck dissection: a multicenter longitudinal study by the Japanese Clinical Study Group on Standardization of Treatment for Lymph Node Metastasis of Head and Neck Cancer, *International Journal of Clinical Oncology*, **15**, 33–38.
- Oakes, D. (1989), Bivariate Survival Models Induced by Frailties, *Journal of American Statistical Association*, **84**, 487–493.
- Parner, E. (1998), Asymptotic Theory for the Correlated Gamma-frailty Model, *The Annals of Statistics*, **26**, 183–214.
- Patterson, H.D. and Thompson, R. (1974), Recovery of Interblock Information When Block Sizes are Unequal, *Biometrika*, **58**, 545–554.

- Pawitan, Y. and Self, S. (1993), Modeling Disease Marker Processes in AIDS, *Journal of American Statistical Association*, **83**, 719–726.
- Pollard, D. (1990), *Empirical Processes: Theory and Applications*, Hayward, CA: Institute of Mathematical Statistics.
- Prentice, R.L. (1982), Covariate Measurement Errors and Parameter Estimation in a Failure Time Regression Model, *Biometrika*, **69**, 331–342.
- Prentice, R.L. and Breslow, N.E. (1978), Retrospective Studies and Failure Time Models, *Biometrika*, **65**, 153–158.
- Ratcliffe, S.J., Guo, W. and Ten Have, T.R. (2004), Joint Modeling of Longitudinal and Survival Data via a Common Frailty, *Biometrics*, **60**, 892–899.
- Ribaudo, H.J., Thompson, S.G., and Allen-Mersh, T.G. (2000), A Joint Analysis of Quality of Life and Survival Using a Random Effect Selection Model, *Statistics in Medicine*, **19**, 3237–3250.
- Ringash, J., Bezjak, A., O’Sullivan, B., and Redelmeier, D.A. (2004), Interpreting Differences in Quality of Life: The FACT-H&N in Laryngeal Cancer Patients, *Quality of Life Research*, **13**, 725–733.
- Ripatti, S. and Palmgren, J. (2000), Estimation of multivariate frailty models using penalized partial likelihood, *Biometrics*, **56**, 1016–1022.
- Rizopoulos, D., Verbeke, G., Lesaffre, E., and Vanrenterghem, Y. (2008), A Two-Part Joint Model for the Analysis of Survival and Longitudinal Binary Data with Excess Zeros, *Biometrics*, **64**, 611–619.
- Rizopoulos, D., Verbeke, G., and Molenberghs, G. (2008), Shared Parameter Models under Random Effects Misspecification, *Biometrika*, **95**, 63–74.
- Robinson, G.K. (1991), That BLUP is a Good Thing: The Estimation of Random Effects, *Statistical Science*, **6**, 15–51.
- Rothman, K.J. (2002), *Epidemiology: An Introduction*, Oxford University Press.
- Schall, R. (1991), Estimation in Generalized Linear Models with Random Effects, *Biometrika*, **78**, 719–727.

- Sen, P.K. and Singer, J.M. (1993), *Large Sample Methods in Statistics*, Chapman & Hall, New York.
- Solomon, P.J. and Cox, D.R. (1992), Nonlinear Components of Variance Models, *Biometrika*, **79**, 1–11.
- Song, X., Davidian, M., and Tsiatis, A.A. (2002), “A Semiparametric Likelihood Approach to Joint Modeling of Longitudinal and Time-to-Event Data, *Biometrics*, **58**, 742–753.
- Song, X. and Wang, C.Y. (2007), Semiparametric Approaches for Joint Modeling of Longitudinal and Survival Data with Time-Varying Coefficients, *Biometrics*, **64**, 557–566.
- Stefanski, L.A. and Carroll, R.J. (1987), Conditional Scores and Optimal Scores for Generalized Linear Measurement- Error Models, *Biometrika*, **74**, 703–716.
- Stiratelli, R., Laird, N., and Ware, J.H. (1984), Random Effects Models for Serial Observations with Binary Responses, *Biometrics*, **40**, 961–971.
- Terrell, J.E., Ronis, D.L., Fowler, K.E., Bradford, C.R., Chepeha, D.B., Prince, M.E., Teknos, T.N., Wolf, G.T., and Duffy, S.A. (2004), Clinical predictors of quality of life in patients with head and neck cancer, *Archives of Otolaryngology Head & Neck Surgery*, **130**, 401–408.
- Therneau, T.M., and Grambsch, P.M. (2001), *Modeling Survival Data*, New York: Springer.
- Thisted, R.A. (1988), *Elements of Statistical Computing*, Chapman & Hall.
- Tierney, L. and Kadane, J.B. (1986), Accurate Approximations for Posterior Moments and Marginal Densities, *Journal of American Statistical Association*, **81**, 82–86.
- Troxel, A.B. and Esserman, D.A. (2004), Frailty Models for Quality of Life in Oncology, *Journal of Biopharmaceutical Statistics*, **14**, 145–154.
- Tseng, Y.K., Hsieh, R., and Wang, J.L. (2005), Joint Modelling of Accelerated Failure Time and Longitudinal Data, *Biometrika*, **92**, 587–603.
- Tsiatis, A.A. (1981), A Large Sample Study of Cox’s Regression Model, *The Annals of*

*Statistics*, **9**, 93–108.

Tsiatis A.A., Degruetola, V., and Wulfsohn M. (1995), Modeling the Relationship of Survival to Longitudinal Data Measured with Error. Applications to Survival and CD4 Counts in Patients with AIDS, *Journal of American Statistical Association*, **90**, 27–37.

Tsiatis, A.A. and Davidian M. (2001), A Semiparametric Estimator for the Proportional Hazards Model with Longitudinal Covariates Measured with Error, *Biometrika*, **88**, 447–458.

van der Vaart, A.W. (1998), *Asymptotic Statistics*. Cambridge University Press.

van der Vaart, A.W. and Wellner, J.A. (1996), *Weak Convergence and Empirical Processes*. New York: Springer-Verlag.

Verbeke, G. and Lesaffre, E. (1996), A Linear Mixed-effects Model with Heterogeneity in the Random-effects Model with Heterogeneity in the Random-effects Population, *Journal of American Statistical Association*, **91**, 217–221.

Verbeke, G. and Molengerghs, G. (2000), *Linear Mixed Models for Longitudinal Data*, Springer Series in Statistics, Springer-Verlag, New-York.

Waclawiw, M.A. and Liang, K.Y. (1993), Prediction of Random Effects in the Generalized Linear Model, *Journal of American Statistical Association*, **88**, 171–178.

Wang, Y. and Taylor, J.M.G. (2001), Jointly Modeling Longitudinal and Event Time Data with Application to Acquired Immunodeficiency Syndrome, *Journal of American Statistical Association*, **96**, 895–905.

Wei, L.J., Lin, D.Y. and Weissfeld, L. (1989), Regression Analysis of Multivariate Incomplete Failure Time Data by Modeling Marginal Distributions, *Journal of American Statistical Association*, **84**, 1065–1073.

Wu, M. and Bailey, K. (1989), Estimation and Comparison of Changes in the Presence of Informative Right Censoring: Conditional Linear Model, *Biometrics*, **45**, 939–955.

Wu, M. and Carroll, R. (1988), Estimation and Comparison of Changes in the Presence of Informative Right Censoring by Modelling the Censoring Process, *Biometrics*,

**44**, 175–188.

Wulfsohn M. and Tsiatis A.A. (1997), A Joint Model for Survival and Longitudinal Data Measured with Error, *Biometrics*, **53**, 330–339.

Xu, J. and Zeger S. (2001a), The Evaluation of Multiple Surrogate Endpoints, *Biometrics*, **57**, 81–87.

———(2001b), “Joint Analysis of Longitudinal Data Comprising Repeated Measures and Times to Events, *Applied Statistics*, **50**, 375–387.

Zeger, S.L. and Karim, M.R. (1991), Generalized Linear Models with Random Effects: a Gibbs Sampling Approach, *Journal of American Statistical Association*, **86**, 79–86.

Zeger, S.L., Liang, K.Y. and Albert, P.S. (1988), Models for Longitudinal Data: A Generalized Estimating Equation Approach, *Biometrics*, **44**, 1049–1060.

Zeng, D. and Cai, J. (2005a), Simultaneous Modelling of Survival and Longitudinal Data with an Application to Repeated Quality of Life Measures, *Lifetime Data Analysis*, **11**, 151–174.

———(2005b), Asymptotic Results for Maximum Likelihood Estimators in Joint Analysis of Repeated Measurements and Survival Time, *The Annals of Statistics*, **33**, 2132–2163.

Zhang, D. and Davidian, M. (2001), Linear Mixed Models with Flexible Distributions of Random Effects for Longitudinal Data, *Biometrics*, **57**, 795–802.

Ye, W., Lin, X.H., Taylor, and J.M.G. (2008), A penalized likelihood approach to joint modeling of longitudinal measurements and time-to-event data, *Statistics and Its Interface*, **1**, 33–45.